

LEAGER PROGRAMMING

Anthony Chyr

UUCS-19-002

School of Computing
University of Utah
Salt Lake City, UT 84112 USA

February 26, 2019

Abstract

Leager programming, a portmanteau of “lazy” and “eager” or “limit” and “eager,” is an evaluation strategy that mixes lazy evaluation and eager evaluation. This evaluation strategy allows iterators to precompute the next value in a separate thread, storing the result in a cache until it is needed by the caller. Leager programming often takes the form of an iterator, which alone allows data to be prefetched, and when chained together can be used to form concurrent pipelines. We found a dramatic reduction in latency on par with code written with asynchronous callbacks, while making minimal modifications to the initial sequential code.

LEAGER PROGRAMMING

by

Anthony Chyr

A senior thesis submitted to the faculty of
The University of Utah
in partial fulfillment of the requirements for the degree of

Bachelor of Computer Science

School of Computing
The University of Utah
May 2019

Approved:

_____/_____
Matthew Flatt
Supervisor

_____/_____
H. James de St. Germain
Director of Undergraduate Studies
School of Computing

_____/_____
Ross Whitaker
Diretor
School of Computing

Copyright © Anthony Chyr 2019

All Rights Reserved

ABSTRACT

Leager programming, a portmanteau of “lazy” and “eager” or “limit” and “eager,” is an evaluation strategy that mixes lazy evaluation and eager evaluation. This evaluation strategy allows iterators to precompute the next value in a separate thread, storing the result in a cache until it is needed by the caller. Leager programming often takes the form of an iterator, which alone allows data to be prefetched, and when chained together can be used to form concurrent pipelines. We found a dramatic reduction in latency on par with code written with asynchronous callbacks, while making minimal modifications to the initial sequential code.

CONTENTS

ABSTRACT	iii
LIST OF CODE LISTINGS	v
LIST OF TABLES	vii
CHAPTERS	
1. INTRODUCTION	1
2. BACKGROUND	3
2.1 Eager evaluation vs lazy evaluation	3
2.2 Generators	4
2.3 The space in between eager evaluation and lazy evaluation	5
3. RELATED WORK	6
3.1 The Toyota Production System	6
3.2 Analogies in hardware and software	7
3.3 Automatic parallelization and parallelism by annotation	7
3.4 Futures (programming construct)	8
3.5 Python <code>concurrent.futures</code>	10
3.6 <code>async</code> and <code>await</code>	11
3.7 PyPI <code>prefetch_generator</code> , <code>async_prefetch</code> , and <code>pythonflow</code>	11
4. METHODS	13
4.1 Python decorators	13
4.2 <code>leager</code> and <code>LeagerIterator</code>	14
4.3 <code>lmap</code>	17
4.4 <code>lmap_unordered</code>	17
5. RESULTS	19
5.1 Example: prefetching data	20
5.2 Example: building a pipeline	28
5.3 Python Global Interpreter Lock (GIL)	36
6. CONCLUSIONS	44
APPENDICES	
A. PYTHON IMPLEMENTATION OF LEAGER PROGRAMMING ...	45
B. HELPER LIBRARIES	50
REFERENCES	52

LIST OF CODE LISTINGS

2.1	A prototypical test used to determine whether the language is eager or lazy. . . .	3
2.2	Creating laziness in an eager language by <i>thunking</i> the expression.	4
2.3	An enumeration of the natural numbers expressed as a generator function and as a generator expression.	4
2.4	Extracting values from generator functions and generator expressions.	4
3.1	Example of parallelism by annotation in OpenMP.	8
3.2	A recursively defined Fibonacci function using futures. Deadlocks at <code>fib(5)</code> . . .	9
4.1	Python decorators.	13
4.2	Python decorators desugared.	13
4.3	Example of a decorator that modifies a function such that it prints how long it runs every time it is called.	14
4.4	Example of a lazily evaluated Python generator.	14
4.5	Example of a lazily evaluated Python generator converted into a leagerly evalu- ated Python generator.	15
4.6	Example of a leagerly evaluated generator with a larger cache.	15
4.7	Example of a leager generator leaving scope and getting garbage collected.	16
4.8	Reverting a leager generator back to being lazily evaluated.	16
4.9	Example use of <code>lmap</code>	17
4.10	Example use of <code>lmap_unordered</code>	18
5.1	Example program retrieving articles on Wikipedia.	21
5.2	Using leager programming to prefetch archived articles on Wikipedia.	22
5.3	Using <code>concurrent.futures</code> to prefetch archived articles on Wikipedia.	23
5.4	Using <code>async</code> and <code>await</code> to prefetch archived articles on Wikipedia.	24
5.5	Using <code>prefetch_generator</code> to prefetch archived articles on Wikipedia.	26
5.6	Using <code>pythonflow</code> to retrieve archived articles on Wikipedia.	27
5.7	Example program downloading archived articles on Wikipedia to local storage. .	29
5.8	Using leager programming to form a concurrent pipeline to download archived articles on Wikipedia to local storage.	30
5.9	Using <code>concurrent.futures</code> to form a concurrent pipeline to download archived articles on Wikipedia to local storage.	31
5.10	Using <code>async</code> and <code>await</code> to form a concurrent pipeline to download archived articles on Wikipedia to local storage.	32
5.11	Using <code>prefetch_generator</code> to form a concurrent pipeline to download archived articles on Wikipedia to local storage.	34
5.12	Using <code>pythonflow</code> to download archived articles on Wikipedia to local storage. .	35
5.13	Example program with a CPU bound task.	37
5.14	Using leager programming to break a CPU bound task into a concurrent pipeline. Due to the Global Interpreter Lock, no improvement is made.	38
5.15	Using <code>concurrent.futures</code> to break a CPU bound task into a concurrent pipeline. Due to the Global Interpreter Lock, no improvement is made.	39

5.16	Using <code>async</code> and <code>await</code> to break a CPU bound task into a concurrent pipeline. Due to the Global Interpreter Lock, no improvement is made.	40
5.17	Using <code>prefetch_generator</code> to break a CPU bound task into a concurrent pipeline. Due to the Global Interpreter Lock, no improvement is made.	42
5.18	Using <code>pythonflow</code> decompose a CPU bound task into a directed acyclic graph of operations. Due to the Global Interpreter Lock, no improvement is made.	43
A.1	<code>leager.py</code>	45
B.1	<code>funchelp.py</code>	50

LIST OF TABLES

- 5.1 Run times in seconds comparing example programs written using different styles. 19

CHAPTER 1

INTRODUCTION

Many programming tasks have sequential implementations that are pleasant to read, but fail to exploit the opportunities for better performance through the concurrency and parallelism inherent to the task. For example, files can be fetched from a server by issuing requests in a sequential loop, but looking ahead and prefetching the request in parallel can speed up the download by hiding the latency. Similarly, a file processing task can be expressed naturally by looping over a list of files with a sequence of steps that load one file, apply a filter, and save the file back to disk — but streaming files through separate parallel processes for those steps can improve performance using the CPU and I/O in parallel.

Programming language designers and implementers have long recognized the potential for performance improvements by parallelizing otherwise sequential operations. Fully automatic parallelization would be ideal [14], but automatic parallelization has so far succeeded only for certain kinds of problems. Instead of fully automatic parallelization, programmers can recast their problem using features such as asynchronous programming [24], futures [23], and parallel for loops [28]. In those cases, the programmer must adopt a slightly different mental model of the computation, but hopefully one that is not too far from the sequential model.

This text presents another mental model for parallel programming with a particular emphasis on staying close to sequential constructs. It focuses on exploiting the opportunities for concurrency and parallelism inherent in the evaluation of expressions by exploring the space in between two well known evaluation strategies: eager evaluation and lazy evaluation — which I call *leager programming*.

Leager programming is about when: when an expression should be evaluated. In eager evaluation, also called “greedy evaluation” or “strict evaluation,” the expression is evaluated as soon as it is bound to a variable. This presents a problem, for example, when iterating

over an expression that enumerates the natural numbers. In eager evaluation, the entire list is evaluated before iteration begins, which would require infinite time and memory. In lazy evaluation, also called “call-by-need,” the expression is evaluated only when it is actually used by another expression. This allows, for example, iteration over an enumeration of the natural numbers without requiring infinite space, and allows the caller to begin consuming the natural numbers without having to wait an infinite amount of time (but enumerating all elements in the iterator would still take an infinite amount of time).

Leager programming seeks to straddle the space in between these two evaluation strategies. It can be thought of as lazy evaluation that eagerly precomputes the next value, or another way to look at it is it limits the eagerness to only the next value ahead of the consuming caller. For this reason, the term, “leager,” is a portmanteau of “lazy” and “eager” or “limits” and “eager.”

The precomputation in leager programming takes place in a separate thread or threads, concurrent to the consuming caller (or callers). Structured in this way, the thread or threads tasked with producing the precomputed values can be thought of as the eager portion of the leager iterator. Likewise, the thread or threads consuming the precomputed values can be thought of as the lazy portion of the a leager iterator, and the degree to which the leager iterator is lazy or eager can be precisely controlled through thread based synchronization primitives such as locks and signals.

To demonstrate a proof of concept of this evaluation strategy and to ground the concepts and terminology, leager programming was implemented in the Python programming language, specifically CPython 3.6.4, as a library of higher order functions and classes whose source code is given in Appendix A on pg. 45. This implementation encapsulates the complexity of managing the thread overhead, and allows the programmer to convert lazy Python generators and eager Python functions (mapped over an iterator) into a leager iterator by passing them to one of the higher order functions or classes in the leager programming library.

CHAPTER 2

BACKGROUND

Leager programming straddles the space between eager evaluation and lazy evaluation. While some texts restrict the definition of eager and lazy to only refer to when the expressions passed into functions are evaluated [12], this text uses a broader definition that can be applied to all expressions, which is similar to how other texts within the Python community apply the concept [13].

2.1 Eager evaluation vs lazy evaluation

Eager evaluation, also called “greedy evaluation” or “strict evaluation,” evaluates an expression as soon as it is bound to a variable. Lazy evaluation, also called “call-by-need,” defers the evaluation of an expression until when it is actually used by another expression. A prototypical example that highlights the difference is to construct a function that does not use its formal argument as shown in Listing 2.1.

```
1 f = lambda x: None
2 f(1/0)
```

Listing 2.1: A prototypical test used to determine whether the language is eager or lazy.

In languages that eagerly evaluate argument expressions, when the function is called with the expression `1/0` as its argument, the expression is immediately evaluated, which results in a divide by zero exception. On the other hand, in languages that lazily evaluates argument expressions (such as in Haskell), no exception is raised; the expression is never evaluated because the function never uses its formal argument. Because argument expressions are eagerly evaluated in Python, Listing 2.1 results in a divide by zero exception.

Even though Python is an eager language, laziness can be introduced by wrapping the expression in a lambda, also called *thunking*. While the lambda itself is eagerly evaluated, the programmer can explicitly control when its contents, the original expression, is evaluated

by calling the function — thus creating laziness as shown in Listing 2.2.

```
1 f = lambda x: None
2 f(lambda: 1/0)
```

Listing 2.2: Creating laziness in an eager language by *thunking* the expression.

2.2 Generators

In Python, generators can be viewed as another kind of explicit evaluation-control construct, similar to *thunks*, where a generator is created eagerly but its body is evaluated on demand (lazily). An example of these generators is given in Listing 2.3, which implements an enumeration of the natural numbers.

```
1 # generator function
2 def natural_numbers():
3     i = 1
4     while True:
5         yield i
6         i += 1
7
8 # generator expression
9 nat_num = (i for i in natural_numbers())
```

Listing 2.3: An enumeration of the natural numbers expressed as a generator function and as a generator expression.

To evaluate the body of a generator, both generator functions (after it has been applied) and generator expressions are iterators, which means they can be used in a `for` loop or by calling the Python built-in function `next` as shown in Listing 2.4.

```
1 for i in natural_numbers():
2     print(i)
3     if i >= 3:
4         break
5 # 1
6 # 2
7 # 3
8
9 next(nat_num) # 1
10 next(nat_num) # 2
11 next(nat_num) # 3
```

Listing 2.4: Extracting values from generator functions and generator expressions.

2.3 The space in between eager evaluation and lazy evaluation

In eager evaluation, expressions are evaluated whether they are needed or not, and values that may not be needed are computed anyway, consuming both time and memory. This can be disastrous, as expressions that represent, for example, an enumeration of the natural numbers consume infinite time and memory if eagerly evaluated. On the other hand, lazy evaluation does eliminate the waste from computing expressions whose values may not be needed, but deferring the evaluation requires book keeping overhead in addition to losing out on the opportunity of computing the value and having it already available when it is needed.

The space in between eager evaluation and lazy evaluation exists in the time between when an expression is bound to a variable and when it is needed. Leager programming seeks to straddle this space. Similar to generators, leager programming lets the programmer explicitly delay the evaluation, but evaluation is more eager than generators, allowing some values to be computed before they are needed. This takes advantage of the opportunity of having a value already available before it is needed, while also being able to express sequences that have infinite size and limiting the waste of computing values that may not be needed.

CHAPTER 3

RELATED WORK

Improving performance by optimizing when work should be done is not novel, and numerous examples exist within computing and outside of it. While these examples implement aspects of leager programming, the goal of this text is to distill these concepts into its own library, independent of specific application, that expands on public libraries in the Python Package Index (PyPI) [20].

3.1 The Toyota Production System

The initial inspiration for leager programming came from the Toyota Production System. The Toyota Production System, also referred to as lean manufacturing, focuses on the “absolute elimination of waste” [15]. It was a response to the manufacturing practices of the early 20th century, which is often referred to as the “push” method of manufacturing that shares many similarities with eager evaluation.

In push manufacturing, parts are produced in large batches with little regard for the capacity at each stage of the manufacturing process, let alone whether the product will be wanted by the consumer in the end. This created a situation where parts piled into large warehouses, which cost factories valuable space and resulted a logistical nightmare of maintaining massive inventories that in some cases never made it to the customer.

Replace push with eager, warehouse with memory, and process with program — and the similarities between manufacturing parts and computing objects become clear. In response to push manufacturing, the Toyota Production System seeks to approach zero inventory, storing only enough raw materials for that day’s production and emphasizing a “pull” system where only what is needed should be produced [15]. Leager programming is the Toyota Production System applied to computing. It seeks to maintain a cache between stages of computing only big enough for the expected demand for data. For this reason, leager programming tips to the side of lazy evaluation where only what is needed should

be computed (as opposed to trying compute as much as possible within the constraints of time and memory).

3.2 Analogies in hardware and software

In digital logic design, a sequential circuit is comprised of latches that pass through a block of combinational logic back into latches [5]. This can be viewed programmatically as a block of memory feeding into a pure function back into a block of memory. Attach a clock to the circuit, and the circuit starts looking like a generator function consuming an iterator. Extend the sequential circuit into a pipeline and allow for prefetching [17], and how to apply leager programming begins to become apparent. Imperatively, leager programming works by layering a cache between function applications in the same way a sequential circuit layers latches between blocks of combinational logic.

Unix pipes redirect the standard output from one program through a cache and into the standard input of another program [27]. Graphics pipelines pass data through a series of computations with memory shared in between stages [1]. Browsers can overcome the effects of network latency by preloading web pages into caches [10]. Similar to how sequential circuits layer latches between blocks of combinational logic, these specific applications across different domains of computing all share a commonality: the layering of functionality and memory that regulates when work is done — which I claim can be parameterized and distilled into a library of its own.

3.3 Automatic parallelization and parallelism by annotation

While hardware already introduces concurrency and parallelism in machine code through instruction level parallelism [17], another approach for introducing concurrency and parallelism to larger programming constructs would be to use an automatic parallelization tool. These tools often use either compile time or run-time techniques that analyze the source code detecting dependencies and identifying sections of code that can be broken into tasks and run concurrently or in parallel. Such techniques, however, have only succeeded for certain kinds of problems. For this reason, most attention in recent years have been in tools where the programmer annotates parallelizable sections of the program in addition to providing other hints to help the automatic parallelization tool parallelize the program [14].

One such tool is OpenMP, which uses parallelism by annotation as shown in Listing 3.1. OpenMP implements multithreading through a fork-join model where starting with a master thread, which is executed sequentially, OpenMP forks a number of slave threads that divide a problem that can be run in parallel before rejoining the master thread [16]. This approach is shared by leager programming, which likewise uses both parallelism by annotation and worker threads forked from a caller thread.

```

1  int main()
2  {
3      const int SIZE = 1000;
4      int a[SIZE];
5
6      #pragma omp parallel for
7      for (int i = 0; i < SIZE; i++)
8      {
9          a[i] = i;
10     }
11
12     return 0;
13 }

```

Listing 3.1: Example of parallelism by annotation in OpenMP.

Where automatic parallelization tools such as OpenMP differ from leager programming is in its evaluation strategy. Automatic parallelization tools seek to introduce parallelism into what would otherwise be sequential code, retaining the evaluation strategy of the underlying code. Leager programming on the other hand seeks to introduce a new evaluation strategy that uses concurrency and parallelism.

3.4 Futures (programming construct)

Futures — also called promises, delays, deferred, and eventuals — are perhaps the programming construct most similar to leager programming. In fact, applying the definition by Prasad and others, “a future or promise can be thought of as a value that will eventually become available” [18], leager programming may be considered a type of future. More importantly, however, a future takes advantage of the time between when a value is needed, and when it can be evaluated. For example, the expression passed as an argument to a function is known (and can begin execution) at the start of a function’s execution, but may not be used until much later in the function’s execution.

Futures have a long history with unique challenges. Evolving from *thunks* (zero argument functions constructed to delay the evaluation of an expression) in argument expressions [18], futures were first implemented in the Multilisp language as an annotation [7]. One challenge facing futures is scheduling when it should be evaluated. Quite often this is done by the operating system, regulated only by limiting the number of futures that can be evaluated at a time to a particular thread pool [23].

Scheduling futures is a difficult problem — one that has been explored since the original paper on futures [2] to as recently as something Kostyukov faced while developing the *finagle* library at Twitter [11]. The problem with futures isn't in situations where the future performs some operation, usually blocking, and then returns. A scheduler, however simple, will still complete such tasks, even if done suboptimally. The problem is scheduling futures that depend on other futures, such as in recursive functions.

Recursive functions, such as quick sort (as shown in an example in Multilisp [7]), have an attractive property where the parallelism grows at each level of the recursion. This is a double edged sword however, and the challenge is two fold. The first is if the implementation does not restrict the number of futures active at any one time, the stack or heap could overflow as the number of futures grows at each level of recursion. The second is if the implementation does restrict the number of futures that can be active at any one time, the program could deadlock as the futures already active depend on the values of additional futures (that cannot be started due the restriction on the number of futures that can be active at any one time) to unblock. This latter problem is highlighted in Listing 3.2.

```

1 from concurrent.futures import ThreadPoolExecutor
2 executor = ThreadPoolExecutor(max_workers=4)
3
4
5 def fib(n):
6     assert n >= 0, 'n cannot be less than zero'
7     if n < 2:
8         return n
9     fib_1 = executor.submit(fib, n-1)
10    fib_2 = executor.submit(fib, n-2)
11    return fib_1.result() + fib_2.result()

```

Listing 3.2: A recursively defined Fibonacci function using futures. Deadlocks at `fib(5)`.

Kostyukov resolved this problem by prioritizing futures along a specific branch of a

recursion tree [11]. This is analogous to resolving the recursion tree using depth first search as opposed to bread first search (which can cause the stack or heap to overflow). However, this scheduling technique, if applied too aggressively, causes any benefit arising from the concurrency that futures provide to evaporate as the recursion tree would be resolved sequentially. In contrast, leager programming side steps this problem by constraining recursion to iteration, in which the restraints on concurrency (which may be understood as a restraint on eagerness) can be better defined.

3.5 Python `concurrent.futures`

Perhaps the closest standard library to leager programming is the `concurrent.futures` module in the Python standard library [23]. In fact, early versions of the leager programming library were built using `concurrent.futures`. Similar to how leager programming precomputes the next value of an iterator in a separate thread, `concurrent.futures` execute functions concurrently in a separate thread or process. Where the two differ is deciding when such functions should be computed.

Python `concurrent.futures` leans towards eagerness, as soon as a thread or process within its thread or process pool becomes available, the function is scheduled to be computed. On the other hand, a lazy future doesn't make sense. If a program waits until the result of a function is needed, the caller would block as the function is computed, eliminating whatever benefits `concurrent.futures` would have had over sequential code.

Leager programming is in fact a form of call-by-future, but it differs from Python `concurrent.futures` by adding a small amount of intelligence in deciding when such functions should be computed. It neither waits until when the result is needed to begin computation, nor does it attempts to compute all the values. It is somewhere in between: eagerly computing until its cache is full and then waiting until the caller consumes a value before it begins computing again.

Another area where Python `concurrent.futures` differs from leager programming is Python `concurrent.futures` has two main methods for scheduling work: `submit`, which takes a function and its arguments; and `map`, which takes a function and an iterable. While `map` in Python `concurrent.futures` works roughly the same way as `lmap` in leager programming, `submit` takes only a single function and produces only a single value. As

discussed in Section 3.4 on pg. 5, for expressions that produce a single value, the degree to which it is eager or lazy depends on the scheduler, and scheduling futures is a difficult problem. For this reason, leager programming uses a higher order function, `leager`, which takes a generator function that produces multiple values instead of a function that produces a single value.

3.6 `async` and `await`

Concurrency is a core feature of leager programming, which for this reason shares similarities with `async` and `await` found in many languages. First appearing in C# 5 in 2012 [3], `async` and `await` has spread to other languages such as JavaScript (ECMA-262) [6] and Python 3.5 [24]. It is a recent revival of an old concept: cooperative multitasking (also called coroutines). In contrast to futures, which achieves concurrency through preemptive multitasking with each task assigned to a separate thread, `async` and `await` achieves concurrency through cooperative multitasking by yielding control to an event loop within a single thread. This allows concurrency without the cost of context switching between threads.

`async` and `await` are implemented in the `asyncio` package in the Python Standard Library [24], and allow tasks to yield control back to an event loop through the keyword `await` for functions tagged with `async`. This is similar to how producer threads and consumer threads in leager programming are scheduled by the operating system to allow for concurrency. Where the two differ is leager programming uses its concurrency to eagerly precompute values into a cache while `asyncio` lazily finds something else to do when it hits a blocked awaitable task. For this reason, where `concurrent.futures` leans eager, `asyncio` leans lazy, and leager programming seeks to be somewhere in between.

3.7 PyPI `prefetch_generator`, `async_prefetch`, and `pythonflow`

Perhaps the closest library to leager programming is the `prefetch_generator` package in the Python Package Index (PyPI) [9] and the `async_prefetch` recipe on Nikki Bowe's blog [4]. Both `prefetch_generator` and `async_prefetch` implement a decorator for generator functions that uses a producer thread to populate a queue with precomputed values. The leager programming library provides the same feature, but in a cleaner implementation that

allows garbage collection if the caller leaves scope before the generator function terminates in addition to more precise control over the number of precomputed values at any one time.

Where the leager programming library differs the most from `prefetch_generator` and `async_prefetch`, however, is that in addition to using eagerness to precompute the values of lazy generators, laziness is used to throttle the values produced by eager functions. This is done through `lmap` and `lmap_unordered` in the leager programming library where it applies an eager Python function over an iterator.

Leager programming also expands on the prefetch use case in `prefetch_generator` and `async_prefetch` to include chaining leager iterators together to form concurrent pipelines. This pushes leager programming in the direction of data flow programming similar to the `pythonflow` package in the Python Package Index (PyPI) [8] — although not to the extreme where the programmer is required to define a directed acyclic graph of operations.

CHAPTER 4

METHODS

The leager programming library is comprised of the following higher order functions and classes intended to convert lazy Python generators and eager Python functions (mapped over an iterator) into leagerly evaluated iterators:

- `leager` and `LeagerIterator`
- `lmap`
- `lmap_unordered`.

The source code for the python implementation of leager programming is given in Appendix A on pg. 45. A helper library was developed to simplify the development of decorator functions for the Python implementation of leager programming library. The source code for this helper library is given in Appendix B on pg. 50.

4.1 Python decorators

For those who may be unfamiliar with Python decorators, this section is intended to be a review. A Python decorator is pure syntactic sugar used to modify the behavior of a function or class [26]. For example, Listing 4.1 is equivalent to Listing 4.2.

```
1 @decorator
2 def decorated():
3     pass
```

Listing 4.1: Python decorators.

```
1 def decorated():
2     pass
3
4 decorated = decorator(decorated)
```

Listing 4.2: Python decorators desugared.

Decorators are defined as high order functions that receive the decorated function as an argument, and composes a new object to be returned as the decorated function's replacement. For example, Listing 4.3 defines a decorator that modifies a function such that it prints how long it runs every time it is called.

```
1 from time import time
2
3
4 def print_runtime(func):
5     def timed_func(*args, **kwargs):
6         start_time = time()
7         return_value = func(*args, **kwargs)
8         print(time() - start_time)
9         return return_value
10    return timed_func
```

Listing 4.3: Example of a decorator that modifies a function such that it prints how long it runs every time it is called.

4.2 leager and LeagerIterator

Python generators are defined as functions with one or more `yield` statements, an example is shown in Listing 4.4. It is lazily evaluated, beginning execution when a value is needed in an iteration and pausing after a value has been yielded. Because Python generators are defined as functions, to modify its behavior, `leager` was defined as a higher order function intended to be used as a Python decorator.

```
1 def example_generator():
2     i = 1
3     while True:
4         yield i
5         i += 1
```

Listing 4.4: Example of a lazily evaluated Python generator.

```

1 from leager import *
2
3
4 @leager
5 def example_generator():
6     i = 1
7     while True:
8         yield i
9         i += 1

```

Listing 4.5: Example of a lazily evaluated Python generator converted into a leagerly evaluated Python generator.

`leager` receives the generator function as one of its arguments, and then composes a leager version of the generator function that it then uses to replace the original generator function. Thus, converting a lazy Python generator function into a leager Python generator function as shown in Listing 4.5. The degree to which a leager generator is eager or lazy can be precisely controlled by adjusting the cache size, which is passed as an optional first argument to the `leager` decorator as shown in Listing 4.6.

```

1 @leager(5)
2 def example_generator():
3     i = 1
4     while True:
5         yield i
6         i += 1

```

Listing 4.6: Example of a leagerly evaluated generator with a larger cache.

Inside the composition, `leager` applies the generator function to acquire its iterator, which it then uses to initialize `LeagerIterator`. `LeagerIterator` maintains the cache, which takes the form of a queue, in addition to starting two daemon threads. The first is the eager portion of the leager iterator, which precomputes values until the cache is full. The second is a custom garbage collector.

If the caller leaves scope before the eager portion of the `LeagerIterator` finishes consuming the iterator as shown in Listing 4.7, thus releasing its reference, the `LeagerIterator` will still have references in two other threads — thus preventing the object from being freed by Python’s garbage collector. In order to rectify this memory leak, the custom garbage collector regularly checks the number of references to the `LeagerIterator`. If the number of references drops to the number of references in the daemon threads, the custom garbage

collector signals for both daemon threads to terminate, thus freeing the references and allowing the `LeagerIterator` to be freed by Python's garbage collector.

```

1 def example_scope():
2     for i in example_generator():
3         print(i)
4         if i >= 3:
5             return # example_generator leaves scope, eager threads terminate
6
7 example_scope()

```

Listing 4.7: Example of a leager generator leaving scope and getting garbage collected.

The two daemon threads are implemented using the `threading` package in Python's standard library [21]. The producing eager thread and the consuming caller thread maintain synchronization through Python `Lock` and `Condition` objects. A shared `Lock` maintains the consistency of the cache, while `Condition` allows the producer and consumer to notify one another when the cache has been mutated. The producer blocks when the cache is full, and is notified by the consumer when it dequeues a value. Likewise, if the cache is empty, the consumer blocks, and the producer notifies the consumer when a value becomes available.

`LeagerIterator` contains a `stop` function that signals the producing eager thread to terminate, thus reverting the behavior of a `LeagerIterator` back to being lazy as shown in Listing 4.8.

```

1 leager_gen = example_generator()
2
3 for i in leager_gen: # leagerly evaluated
4     print(i)
5     if i >= 3:
6         leager_gen.stop() # stop eager threads, revert behavior back to lazy
7         break
8
9 for i in leager_gen: # lazily evaluated
10    print(i)
11    if i >= 6:
12        break

```

Listing 4.8: Reverting a leager generator back to being lazily evaluated.

While `leager` adds eagerness to inherently lazy generators to produce a leager iterator, laziness can be added to inherently eager functions (applied over an iterator) to produce

a leager iterator — which is the direction `lmap` and `lmap_unordered` approaches leager programming.

4.3 `lmap`

Similar to Python’s built-in `map`, which is lazily evaluated [25], and `imap` in the `multi-processing` module of Python’s standard library, which is eagerly evaluated [22] — `lmap` leagerly applies a function over an iterator as shown in Listing 4.9. Similar to `leager`, the degree to which `lmap` is eager or lazy can be precisely controlled by adjusting the size of the cache size, which is passed as an optional third argument.

```

1 from random import randint
2 from time import sleep
3 from leager import *
4
5
6 def square(i):
7     sleep(randint(0, 5))
8     return i * i
9
10 for sqr in lmap(square, range(20), 5):
11     print(sqr)
12     input('Press [enter] to show the next perfect square.')
```

Listing 4.9: Example use of `lmap`.

Unlike `LeagerIterator`, `lmap` does not use a single eager producer thread, but instead spawns a thread for each function application. Each function application is given an index, which is used to ensure that the order in which values are yielded matches that of the iterator `lmap` is mapped over. Because each value is associated with an index, the cache uses a dictionary instead of a queue. And because each function application takes place in its own thread, a higher degree of concurrency can be achieved when compared to `leager`.

Synchronization is maintained through `Lock` and `Condition` objects, similar to `LeagerIterator`. However, since there is no single producer thread to notify, the consumer thread spawns producer threads as needed in order to maintain the cache.

4.4 `lmap_unordered`

Similar to `imap_unordered` in the `multiprocessing` module of Python’s standard library [22], which is eagerly evaluated — `lmap_unordered` leagerly applies a function over

an iterator. Similar to `lmap`, the degree to which `lmap_unordered` is eager or lazy can be precisely controlled by adjusting the size of the cache, which is passed as an optional third argument as shown in Listing 4.10.

```
1 from random import randint
2 from time import sleep
3 from leager import *
4
5
6 def square(i):
7     sleep(randint(0, 5))
8     return i * i
9
10 for sqr in lmap_unordered(square, range(20), 5): # may come out of order
11     print(sqr)
12     input('Press [enter] to show another perfect square.')
```

Listing 4.10: Example use of `lmap_unordered`.

Unlike `lmap`, order is not maintained. No function application is given an index, and adds its value to the cache as soon as it becomes available. For this reason, `lmap_unordered` can be faster than `lmap`. Because order is not maintained, the cache uses a queue instead of a dictionary similar to `LeagerIterator`.

Synchronization is maintained similar to `lmap`.

CHAPTER 5

RESULTS

A summary of the run times comparing the different programming styles for the examples discussed in this chapter is given in Table 5.1.

Style	Prefetch example ¹	Pipeline example ²	CPU-bound example ³
Regular Python	0.67	6.75	2.16
Leager programming	0.12	0.86	2.16
<code>concurrent.futures</code>	0.12	0.85	2.16
<code>async</code> and <code>await</code>	0.01	0.88	2.16
<code>prefetch_generator</code>	0.12	5.63	2.16
<code>pythonflow</code>	0.67	6.76	2.16

Table 5.1: Run times in seconds comparing example programs written using different styles.

¹See Section 5.1 on pg. 20.

²See Section 5.2 on pg. 28.

³See Section 5.3 on pg. 36.

5.1 Example: prefetching data

Leager programming can be used to prefetch data, which may dramatically improve the performance of programs with blocking operations. An example program is shown in Listing 5.1, which retrieves and prints archived articles on Wikipedia with random prime ids. This program is then accelerated using leager programming to prefetch articles in Listing 5.2, which results in a decrease in latency from 0.67 seconds to 0.12 seconds.

As shown in Listing 5.3, this decrease in latency is similarly observed when using `concurrent.futures` to accelerate the program, taking 0.12 seconds. However, `concurrent.futures` requires a more extensive rewrite of the sequential code in Listing 5.1 when compared to leager programming. Using `async` and `await` as shown in Listing 5.4 resulted in the lowest latency of 0.01 seconds by avoiding thread overhead, but requires the entire program to be rewritten using coroutines. As shown in Listing 5.5, `prefetch_generator` does not come with a `map`, but one can easily be created by composing `map` in Python with `background` in `prefetch_generator`. Using `prefetch_generator` resulted in the same latency of 0.12 seconds as leager programming. `pythonflow` is not a library aimed at improving performance or creating concurrency. It instead focuses on introducing dataflow programming to Python. For this reason, `pythonflow`, as shown in Listing 5.6 and with a latency of 0.67, did not result in any improvement when compared to the regular python for this example.

```
1 from random import randint
2 from urllib.request import urlopen
3 from time import time
4
5 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
6
7
8 def is_prime(num):
9     for i in range(2, num):
10         if not num % i:
11             return False
12     return True
13
14
15 def random_prime():
16     while True:
17         num = randint(0, 1000000)
18         if is_prime(num):
19             yield num
20
21
22 def main():
23     start_time = time()
24     for req in map(urlopen, (base_url % oldid for oldid in random_prime())):
25         print(req.read())
26         print('Execution time %s seconds' % (time() - start_time))
27         input('Press [enter] to show the next article.\n')
28         start_time = time()
29
30
31 if __name__ == '__main__':
32     main()
```

Listing 5.1: Example program retrieving articles on Wikipedia.

```

1 from random import randint
2 from urllib.request import urlopen
3 from time import time
4 from leager import *
5
6 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
7
8
9 def is_prime(num):
10     for i in range(2, num):
11         if not num % i:
12             return False
13     return True
14
15
16 @leager
17 def random_prime():
18     while True:
19         num = randint(0, 1000000)
20         if is_prime(num):
21             yield num
22
23
24 def main():
25     start_time = time()
26     for req in lmap(urlopen, (base_url % oldid for oldid in random_prime())):
27         print(req.read())
28         print('Execution time %s seconds' % (time() - start_time))
29         input('Press [enter] to show the next article.\n')
30         start_time = time()
31
32
33 if __name__ == '__main__':
34     main()

```

Listing 5.2: Using leager programming to prefetch archived articles on Wikipedia.

```

1 from random import randint
2 from urllib.request import urlopen
3 from time import time
4 from concurrent.futures import ThreadPoolExecutor
5
6 pool = ThreadPoolExecutor()
7 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
8
9
10 def is_prime(num):
11     for i in range(2, num):
12         if not num % i:
13             return False
14     return True
15
16
17 def random_prime():
18     next_num = pool.submit(randint, 0, 1000000)
19     while True:
20         num = next_num.result()
21         next_num = pool.submit(randint, 0, 1000000)
22         if is_prime(num):
23             yield num
24
25
26 def main():
27     random_prime_generator = random_prime()
28     next_req = pool.submit(urlopen, base_url % next(random_prime_generator))
29     while True:
30         start_time = time()
31         req = next_req.result()
32         next_req = pool.submit(urlopen, base_url % next(random_prime_generator))
33         print(req.read())
34         print('Execution time %s seconds' % (time() - start_time))
35         input('Press [enter] to show the next article.\n')
36
37
38 if __name__ == '__main__':
39     main()

```

Listing 5.3: Using `concurrent.futures` to prefetch archived articles on Wikipedia.

```
1 from random import randint
2 from urllib.request import urlopen
3 from time import time
4 import asyncio
5
6 oldid_queue = asyncio.Queue(1)
7 request_queue = asyncio.Queue(1)
8 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
9
10
11 def is_prime(num):
12     for i in range(2, num):
13         if not num % i:
14             return False
15     return True
16
17
18 def random_prime():
19     while True:
20         num = randint(0, 1000000)
21         if is_prime(num):
22             yield num
23
24
25 async def oldid_producer(loop):
26     rp = random_prime()
27     while True:
28         oldid = await loop.run_in_executor(None, lambda: next(rp))
29         await oldid_queue.put(oldid)
30
31
32 async def request_producer(loop):
33     while True:
34         oldid = await oldid_queue.get()
35         request = await loop.run_in_executor(None, urlopen, base_url % oldid)
36         await request_queue.put(request)
37
38
39 async def user_prompt(loop):
40     while True:
41         start_time = time()
42         request = await request_queue.get()
43         print(request.read())
44         print('Execution time %s seconds' % (time() - start_time))
45         print('Press [enter] to show the next article.')
46         await loop.run_in_executor(None, input)
47
48
49 def main():
50     loop = asyncio.get_event_loop()
51     loop.create_task(oldid_producer(loop))
52     loop.create_task(request_producer(loop))
53     loop.create_task(user_prompt(loop))
```



```
54     loop.run_forever()
55     loop.close()
56
57
58 if __name__ == '__main__':
59     main()
```

Listing 5.4: Using `async` and `await` to prefetch archived articles on Wikipedia.

```

1 from random import randint
2 from urllib.request import urlopen
3 from time import time
4 from prefetch_generator import background
5
6 pgmap = lambda f, it, max_size=1: background(max_size)(map)(f, it)
7 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
8
9 def is_prime(num):
10     for i in range(2, num):
11         if not num % i:
12             return False
13     return True
14
15
16 @background()
17 def random_prime():
18     while True:
19         num = randint(0, 1000000)
20         if is_prime(num):
21             yield num
22
23
24 def main():
25     start_time = time()
26     for req in pgmap(urlopen, (base_url % oldid for oldid in random_prime())):
27         print(req.read())
28         print('Execution time %s seconds' % (time() - start_time))
29         input('Press [enter] to show the next article.\n')
30         start_time = time()
31
32
33 if __name__ == '__main__':
34     main()

```

Listing 5.5: Using `prefetch_generator` to prefetch archived articles on Wikipedia.

```

1 from random import randint
2 from urllib.request import urlopen
3 from time import time
4 import pythonflow as pf
5
6 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
7
8
9 def is_prime(num):
10     for i in range(2, num):
11         if not num % i:
12             return False
13     return True
14
15
16 def random_prime():
17     while True:
18         num = randint(0, 1000000)
19         if is_prime(num):
20             yield num
21
22
23 def main():
24     with pf.Graph() as graph:
25         oldid = pf.placeholder(name='oldid')
26         url = pf.map_(lambda oldid: base_url % oldid, oldid)
27         request = pf.map_(urlopen, url)
28
29         start_time = time()
30         for req in graph(request, oldid=random_prime()):
31             print(req.read())
32             print('Execution time %s seconds' % (time() - start_time))
33             input('Press [enter] to show the next article.\n')
34             start_time = time()
35
36
37 if __name__ == '__main__':
38     main()

```

Listing 5.6: Using pythonflow to retrieve archived articles on Wikipedia.

5.2 Example: building a pipeline

Leager programming can be used to form concurrent pipelines. An example program is shown in Listing 5.7 as archived articles are downloaded from Wikipedia and saved to local storage. This program is then accelerated using leager programming to form concurrent pipelines allowing for a high degree of both horizontal and vertical parallelism as shown in Listing 5.8. This results in a decrease in run time from 6.75 seconds to 0.86 seconds.

As shown in Listing 5.9, this decrease in run time is similarly observed when using `concurrent.futures`, taking 0.85 seconds, but with a more extensive rewrite of the sequential code when compared to leager programming in Listing 5.8. The decrease in runtime is also shared when using `async` and `await` as shown in Listing 5.10, taking 0.88 seconds, but it resulted in the entire program being rewritten using coroutines. `prefetch_generator`, as shown in Listing 5.11, did result in a decrease in run time taking 5.63 seconds, but the high level of horizontal concurrency found in leager programming could not be created by composing `background` in `prefetch_generator` with `map` in Python. `pythonflow`, as shown in Listing 5.12, is not a performance library and did not result in an improvement in performance when compared to regular python, taking 6.76 seconds.

```
1 from urllib.request import urlopen
2 from time import time
3
4 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
5
6
7 def make_request(oidid):
8     req = urlopen(base_url % oidid)
9     req.oidid = oidid
10    return req
11
12
13 def save_request(req):
14     fn = 'scratch/%s.htm' % req.oidid
15     with open(fn, 'wb') as f:
16         return 'fn: %s, bytes_written: %s' % (fn, f.write(req.read()))
17
18
19 def main():
20     start_time = time()
21
22     stage_1 = map(make_request, range(16))
23     stage_2 = map(save_request, stage_1)
24
25     for status_message in stage_2:
26         print(status_message)
27
28     print('Execution time: %s seconds' % (time() - start_time))
29
30
31 if __name__ == '__main__':
32     main()
```

Listing 5.7: Example program downloading archived articles on Wikipedia to local storage.

```
1 from urllib.request import urlopen
2 from time import time
3 from leager import *
4
5 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
6
7
8 def make_request(oidid):
9     req = urlopen(base_url % oidid)
10    req.oidid = oidid
11    return req
12
13
14 def save_request(req):
15     fn = 'scratch/%s.htm' % req.oidid
16     with open(fn, 'wb') as f:
17         return 'fn: %s, bytes_written: %s' % (fn, f.write(req.read()))
18
19
20 def main():
21     start_time = time()
22
23     stage_1 = lmap_unordered(make_request, range(16), 16)
24     stage_2 = lmap_unordered(save_request, stage_1, 16)
25
26     for status_message in stage_2:
27         print(status_message)
28
29     print('Execution time: %s seconds' % (time() - start_time))
30
31
32 if __name__ == '__main__':
33     main()
```

Listing 5.8: Using leager programming to form a concurrent pipeline to download archived articles on Wikipedia to local storage.

```
1 from urllib.request import urlopen
2 from time import time
3 from concurrent.futures import ThreadPoolExecutor
4
5 pool = ThreadPoolExecutor()
6 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
7
8
9 def make_request(oidid):
10     req = urlopen(base_url % oidid)
11     req.oidid = oidid
12     return req
13
14
15 def save_request(req):
16     fn = 'scratch/%s.htm' % req.oidid
17     with open(fn, 'wb') as f:
18         return 'fn: %s, bytes_written: %s' % (fn, f.write(req.read()))
19
20
21 def main():
22     start_time = time()
23
24     stage_1 = pool.map(make_request, range(16))
25     stage_2 = pool.map(save_request, stage_1)
26
27     for status_message in stage_2:
28         print(status_message)
29
30     print('Execution time: %s seconds' % (time() - start_time))
31
32
33 if __name__ == '__main__':
34     main()
```

Listing 5.9: Using `concurrent.futures` to form a concurrent pipeline to download archived articles on Wikipedia to local storage.

```
1 from urllib.request import urlopen
2 from time import time
3 import asyncio
4
5 OBJ_COUNT = 16
6 stage_1_queue = asyncio.Queue(OBJ_COUNT)
7 stage_2_queue = asyncio.Queue(OBJ_COUNT)
8 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
9
10
11 def make_request(oidid):
12     req = urlopen(base_url % oidid)
13     req.oidid = oidid
14     return req
15
16
17 def save_request(req):
18     fn = 'scratch/%s.htm' % req.oidid
19     with open(fn, 'wb') as f:
20         return 'fn: %s, bytes_written: %s' % (fn, f.write(req.read()))
21
22
23 async def stage(loop, func, args, out_queue):
24     result = await loop.run_in_executor(None, func, *args)
25     await out_queue.put(result)
26
27
28 async def stage_1(loop):
29     for oidid in range(OBJ_COUNT):
30         loop.create_task(stage(loop, make_request, (oidid,), stage_1_queue))
31
32
33 async def stage_2(loop):
34     items_processed = 0
35     while items_processed < OBJ_COUNT:
36         request = await stage_1_queue.get()
37         loop.create_task(stage(loop, save_request, (request,), stage_2_queue))
38         items_processed += 1
39
40
41 async def consumer(loop):
42     items_processed = 0
43     while items_processed < OBJ_COUNT:
44         status_message = await stage_2_queue.get()
45         await loop.run_in_executor(None, print, status_message)
46         items_processed += 1
47     loop.stop()
48
49
50 def main():
51     start_time = time()
52
53     loop = asyncio.get_event_loop()
```



```
54     loop.create_task(stage_1(loop))
55     loop.create_task(stage_2(loop))
56     loop.create_task(consumer(loop))
57     loop.run_forever()
58     loop.close()
59
60     print('Execution time: %s seconds' % (time() - start_time))
61
62
63 if __name__ == '__main__':
64     main()
```

Listing 5.10: Using `async` and `await` to form a concurrent pipeline to download archived articles on Wikipedia to local storage.

```

1 from urllib.request import urlopen
2 from time import time
3 from prefetch_generator import background
4
5 pmap = lambda f, it, max_size=1: background(max_size)(map)(f, it)
6 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
7
8
9 def make_request(oidid):
10     req = urlopen(base_url % oidid)
11     req.oidid = oidid
12     return req
13
14
15 def save_request(req):
16     fn = 'scratch/%s.htm' % req.oidid
17     with open(fn, 'wb') as f:
18         return 'fn: %s, bytes_written: %s' % (fn, f.write(req.read()))
19
20
21 def main():
22     start_time = time()
23
24     stage_1 = pmap(make_request, range(16), 16)
25     stage_2 = pmap(save_request, stage_1, 16)
26
27     for status_message in stage_2:
28         print(status_message)
29
30     print('Execution time: %s seconds' % (time() - start_time))
31
32
33 if __name__ == '__main__':
34     main()

```

Listing 5.11: Using `prefetch_generator` to form a concurrent pipeline to download archived articles on Wikipedia to local storage.

```

1 from urllib.request import urlopen
2 from time import time
3 import pythonflow as pf
4
5 base_url = 'https://en.wikipedia.org/w/index.php?oldid=%s'
6
7
8 def make_request(oldid):
9     req = urlopen(base_url % oldid)
10    req.oldid = oldid
11    return req
12
13
14 def save_request(req):
15     fn = 'scratch/%s.htm' % req.oldid
16     with open(fn, 'wb') as f:
17         return 'fn: %s, bytes_written: %s' % (fn, f.write(req.read()))
18
19
20 def main():
21     with pf.Graph() as graph:
22         oldid = pf.placeholder(name='oldid')
23         stage_1 = pf.map_(make_request, oldid)
24         stage_2 = pf.map_(save_request, stage_1)
25
26     start_time = time()
27
28     for status_message in graph(stage_2, oldid=range(16)):
29         print(status_message)
30
31     print('Execution time: %s seconds' % (time() - start_time))
32
33
34 if __name__ == '__main__':
35     main()

```

Listing 5.12: Using pythonflow to download archived articles on Wikipedia to local storage.

5.3 Python Global Interpreter Lock (GIL)

An important limitation in the Python implementation of leager programming is the Python Global Interpreter Lock (GIL), which allows only one thread to execute in CPython 3.6.4 at a time [19]. This prevents any improvement for CPU bound tasks to be made as both the sequential implementation in Listing 5.13 and the leager implementation execute in about 2.16 seconds.

For the same reason, `concurrent.futures` as shown in Listing 5.15, `async` and `await` as shown in Listing 5.16, `prefetch_generator` as shown in Listing 5.17, and `pythonflow` as shown in Listing 5.18 — all resulted in the same run time as the sequential implementation, taking about 2.16 seconds.

```
1 from time import time
2
3
4 def cpu_bound_task(i):
5     n = 1000000
6     while n > 0:
7         n -= 1
8     return i
9
10
11 def main():
12     start_time = time()
13
14     stage_1 = map(cpu_bound_task, range(10))
15     stage_2 = map(cpu_bound_task, stage_1)
16
17     for i in stage_2:
18         print(i)
19
20     print('Execution time: %s seconds' % (time() - start_time))
21
22
23 if __name__ == '__main__':
24     main()
```

Listing 5.13: Example program with a CPU bound task.

```
1 from time import time
2 from leager import *
3
4
5 def cpu_bound_task(i):
6     n = 1000000
7     while n > 0:
8         n -= 1
9     return i
10
11
12 def main():
13     start_time = time()
14
15     stage_1 = lmap_unordered(cpu_bound_task, range(10), 4)
16     stage_2 = lmap_unordered(cpu_bound_task, stage_1, 4)
17
18     for i in stage_2:
19         print(i)
20
21     print('Execution time: %s seconds' % (time() - start_time))
22
23
24 if __name__ == '__main__':
25     main()
```

Listing 5.14: Using leager programming to break a CPU bound task into a concurrent pipeline. Due to the Global Interpreter Lock, no improvement is made.

```
1 from time import time
2 from concurrent.futures import ThreadPoolExecutor
3
4 pool = ThreadPoolExecutor()
5
6
7 def cpu_bound_task(i):
8     n = 1000000
9     while n > 0:
10         n -= 1
11     return i
12
13
14 def main():
15     start_time = time()
16
17     stage_1 = pool.map(cpu_bound_task, range(10))
18     stage_2 = pool.map(cpu_bound_task, stage_1)
19
20     for i in stage_2:
21         print(i)
22
23     print('Execution time: %s seconds' % (time() - start_time))
24
25
26 if __name__ == '__main__':
27     main()
```

Listing 5.15: Using `concurrent.futures` to break a CPU bound task into a concurrent pipeline. Due to the Global Interpreter Lock, no improvement is made.

```
1 from time import time
2 import asyncio
3
4 OBJ_COUNT = 10
5 stage_1_queue = asyncio.Queue(OBJ_COUNT)
6 stage_2_queue = asyncio.Queue(OBJ_COUNT)
7
8
9 def cpu_bound_task(i):
10     n = 1000000
11     while n > 0:
12         n -= 1
13     return i
14
15
16 async def stage(loop, func, args, out_queue):
17     result = await loop.run_in_executor(None, func, *args)
18     await out_queue.put(result)
19
20
21 async def stage_1(loop):
22     for i in range(OBJ_COUNT):
23         loop.create_task(stage(loop, cpu_bound_task, (i,), stage_1_queue))
24
25
26 async def stage_2(loop):
27     items_processed = 0
28     while items_processed < OBJ_COUNT:
29         result = await stage_1_queue.get()
30         loop.create_task(stage(loop, cpu_bound_task, (result,), stage_2_queue))
31         items_processed += 1
32
33
34 async def consumer(loop):
35     items_processed = 0
36     while items_processed < OBJ_COUNT:
37         result = await stage_2_queue.get()
38         await loop.run_in_executor(None, print, result)
39         items_processed += 1
40     loop.stop()
41
42
43 def main():
44     start_time = time()
45
46     loop = asyncio.get_event_loop()
47     loop.create_task(stage_1(loop))
48     loop.create_task(stage_2(loop))
49     loop.create_task(consumer(loop))
50     loop.run_forever()
51     loop.close()
52
53     print('Execution time: %s seconds' % (time() - start_time))
```



```
54 |
55 |
56 | if __name__ == '__main__':
57 |     main()
```

Listing 5.16: Using `async` and `await` to break a CPU bound task into a concurrent pipeline. Due to the Global Interpreter Lock, no improvement is made.

```
1 from time import time
2 from prefetch_generator import background
3
4 pimap = lambda f, it, max_size=1: background(max_size)(map)(f, it)
5
6
7 def cpu_bound_task(i):
8     n = 1000000
9     while n > 0:
10         n -= 1
11     return i
12
13
14 def main():
15     start_time = time()
16
17     stage_1 = pimap(cpu_bound_task, range(10), 4)
18     stage_2 = pimap(cpu_bound_task, stage_1, 4)
19
20     for i in stage_2:
21         print(i)
22
23     print('Execution time: %s seconds' % (time() - start_time))
24
25
26 if __name__ == '__main__':
27     main()
```

Listing 5.17: Using `prefetch_generator` to break a CPU bound task into a concurrent pipeline. Due to the Global Interpreter Lock, no improvement is made.

```
1 from time import time
2 import pythonflow as pf
3
4
5 def cpu_bound_task(i):
6     n = 1000000
7     while n > 0:
8         n -= 1
9     return i
10
11
12 def main():
13     with pf.Graph() as graph:
14         it = pf.placeholder(name='it')
15         stage_1 = pf.map_(cpu_bound_task, it)
16         stage_2 = pf.map_(cpu_bound_task, stage_1)
17
18         start_time = time()
19
20         for result in graph(stage_2, it=range(10)):
21             print(result)
22
23         print('Execution time: %s seconds' % (time() - start_time))
24
25
26 if __name__ == '__main__':
27     main()
```

Listing 5.18: Using pythonflow decompose a CPU bound task into a directed acyclic graph of operations. Due to the Global Interpreter Lock, no improvement is made.

CHAPTER 6

CONCLUSIONS

Leager programming is about: when an expression should be evaluated. It straddles the space in between eager evaluation and lazy evaluation, improving performance by introducing concurrency into what would otherwise be sequential code.

While improving performance by optimizing when work should be done is not novel, the goal of the leager programming library is to distill these concepts into its own library, independent of specific applications, that expand on publicly available libraries in the Python Package Index (PyPI).

Built using the `threading` package in Python's standard library, it is comprised of higher order functions and classes taking the form of decorators and substitutions of well known functions in Python. Used in asynchronous callbacks, leager programming can be used to prefetch data; and chained together, it can be used to form concurrent pipelines.

APPENDIX A

PYTHON IMPLEMENTATION OF LEAGER PROGRAMMING

```
1  """Tools that intelligently combine eager evaluation and lazy evaluation."""
2
3  from collections import deque, abc
4  from functools import wraps
5  from sys import getrefcount
6  from threading import Lock, Condition, Thread, current_thread
7  from time import sleep
8
9  from funchelp import default, unwrap
10
11
12  __all__ = ['leager', 'LeagerIterator', 'lmap', 'lmap_unordered']
13
14
15  @default
16  @unwrap
17  def leager(gen_func, max_size=1):
18      """
19      Allow generator function to precompute the next value in a separate thread.
20
21      >>> @leager
22      ... def example_generator():
23      ...     from time import sleep
24      ...     while True:
25      ...         sleep(5) # example blocking operation
26      ...         yield
27
28      :param gen_func: generator function
29      :param max_size: maximum number of values to precompute
30      :return: generator function that precomputes
31      """
32      @wraps(gen_func)
33      def leager_gen_func(*args, **kwargs):
34          return LeagerIterator(gen_func(*args, **kwargs), max_size)
35      return leager_gen_func
36
37
38  class LeagerIterator(abc.Iterator):
39      """
40      Create a leager iterator object that precomputes the next value from a
41      provided iterator object in a separate thread.
```

```
42     """
43     def __init__(self, iterable, max_size=1):
44         self.DAEMON_REF_COUNT = 5 # number of references in daemon threads
45         self.GC_CHECK_INTERVAL = 1 # seconds between garbage collection checks
46
47         assert max_size > 0, 'invalid max_size'
48         self.iterator = iter(iterable)
49         self.max_size = max_size
50
51         self.cache = deque()
52
53         self.stop_signal = False
54         self.worker_stopped = False
55
56         self.mutex = Lock()
57         self.consumer = Condition(self.mutex)
58         self.producer = Condition(self.mutex)
59
60         self.worker = Thread(target=self._worker, daemon=True)
61         self.worker.start()
62
63         self.gc = Thread(target=self._gc, daemon=True)
64         self.gc.start()
65
66     def stop(self):
67         """
68         Signal all worker threads to stop precomputing values.
69
70         :return: None
71         """
72         with self.mutex:
73             self._signal_stop()
74
75     def _signal_stop(self):
76         self.stop_signal = True
77         self.producer.notify_all()
78
79     def _stop_worker(self):
80         self.worker_stopped = True
81         self.consumer.notify_all()
82
83     def _worker(self):
84         while True:
85             with self.producer:
86                 if self.stop_signal:
87                     self._stop_worker()
88                     return
89                 while len(self.cache) >= self.max_size:
90                     self.producer.wait()
91                 if self.stop_signal:
92                     self._stop_worker()
93                 return
94
95         try:
```

```

96         item = next(self.iterator)
97     except StopIteration:
98         with self.producer:
99             self._stop_worker()
100         return
101
102     with self.producer:
103         self.cache.append(item)
104         self.consumer.notify()
105
106     def _gc(self):
107         while True:
108             sleep(self.GC_CHECK_INTERVAL)
109             with self.mutex:
110                 if self.worker_stopped:
111                     return
112                 if getrefcount(self) <= self.DAEMON_REF_COUNT:
113                     self._signal_stop()
114
115     def __iter__(self):
116         return self
117
118     def __next__(self):
119         with self.consumer:
120             while not len(self.cache):
121                 if self.worker_stopped:
122                     return next(self.iterator)
123                 self.consumer.wait()
124
125                 self.producer.notify()
126                 return self.cache.popleft()
127
128
129 class lmap(abc.Iterator):
130     """
131     Create a leager map that precomputes the next value in a separate thread.
132     The order in which values are returned is preserved.
133     """
134     class _Counter:
135         def __init__(self):
136             self.count = 0
137
138         def __call__(self):
139             self.count += 1
140             return self.count
141
142     def __init__(self, func, iterable, max_size=1):
143         assert max_size > 0, 'invalid max_size'
144         self.func = func
145         self.iterator = iter(iterable)
146         self.max_size = max_size
147
148         self.produced = self._Counter()
149         self.consumed = self._Counter()

```

```

150
151     self.cache = {}
152     self.workers = set()
153
154     self.mutex = Lock()
155     self.consumer = Condition(self.mutex)
156
157     with self.mutex:
158         self._spawn_workers()
159
160 def _spawn_workers(self):
161     while len(self.cache) + len(self.workers) < self.max_size:
162         try:
163             item = next(self.iterator)
164         except StopIteration:
165             if not len(self.workers):
166                 self.consumer.notify_all()
167             return
168
169         idx = self.produced()
170         worker = Thread(target=self._worker, args=(item, idx), daemon=True)
171         self.workers.add(worker)
172         worker.start()
173
174 def _worker(self, item, idx):
175     value = self.func(item)
176     with self.mutex:
177         self.cache[idx] = value
178         self.workers.remove(current_thread())
179         self.consumer.notify_all()
180
181 def __iter__(self):
182     return self
183
184 def __next__(self):
185     with self.consumer:
186         idx = self.consumed()
187         while idx not in self.cache:
188             if not self.cache and not self.workers:
189                 return self.func(next(self.iterator))
190             self.consumer.wait()
191
192     return_value = self.cache.pop(idx)
193     self._spawn_workers()
194     return return_value
195
196
197 class lmap_unordered(abc.Iterator):
198     """
199     Create a leager map that precomputes the next value in a separate thread.
200     The order in which values are returned is arbitrary.
201     """
202     def __init__(self, func, iterable, max_size=1):
203         assert max_size > 0, 'invalid max_size'

```



```
204     self.func = func
205     self.iterator = iter(iterable)
206     self.max_size = max_size
207
208     self.cache = deque()
209     self.workers = set()
210
211     self.mutex = Lock()
212     self.consumer = Condition(self.mutex)
213
214     with self.mutex:
215         self._spawn_workers()
216
217 def _spawn_workers(self):
218     while len(self.cache) + len(self.workers) < self.max_size:
219         try:
220             item = next(self.iterator)
221         except StopIteration:
222             if not len(self.workers):
223                 self.consumer.notify_all()
224             return
225
226         worker = Thread(target=self._worker, args=(item,), daemon=True)
227         self.workers.add(worker)
228         worker.start()
229
230 def _worker(self, item):
231     value = self.func(item)
232     with self.mutex:
233         self.cache.append(value)
234         self.workers.remove(current_thread())
235         self.consumer.notify()
236
237 def __iter__(self):
238     return self
239
240 def __next__(self):
241     with self.consumer:
242         while not len(self.cache):
243             if not self.workers:
244                 return self.func(next(self.iterator))
245             self.consumer.wait()
246
247     return_value = self.cache.popleft()
248     self._spawn_workers()
249     return return_value
```

Listing A.1: leager.py

APPENDIX B

HELPER LIBRARIES

```
1 from functools import wraps
2
3 __all__ = ['default', 'unwrap']
4
5
6 def default(dec_func):
7     """
8     Allows a decorator function with parameters where all the parameters are
9     optional to use the unapplied form. For example, by specifying
10
11     >>> @default
12     ... def decorator(*args, **kwargs):
13     ... def real_decorator(func):
14     ... @wraps(func)
15     ... def wrapped(*a, **kw):
16     ... func(a, kw, args, kwargs)
17     ... return wrapped
18     ... return real_decorator
19
20     allows
21
22     >>> @decorator
23     ... def decorated():
24     ... pass
25
26     to be equivalent to
27
28     >>> @decorator()
29     ... def decorated():
30     ... pass
31
32     Note that this creates a situation where if the first and only parameter is
33     callable, the behavior is undefined. For example, avoid
34
35     >>> f = lambda *x: x
36     >>> @decorator(f)
37     ... def decorate():
38     ... pass
39
40     :param dec_func: decorator function
41     :return: decorator function that will worked in the unapplied form
42     """
43     @wraps(dec_func)
```

```

44     def default_dec(*args, **kwargs):
45         return dec_func()(args[0]) \
46             if len(args) == 1 and callable(args[0]) and not kwargs else \
47             dec_func(*args, **kwargs)
48     return default_dec
49
50
51 def unwrap(dec_func):
52     """
53     Unwraps a decorator function with parameters forcing the decorated function
54     to be passed as the first argument followed by the decorator's parameters.
55
56     For example, it allows
57
58     >>> def decorator(*args, **kwargs):
59     ... def real_decorator(func):
60     ... @wraps(func)
61     ... def wrapper(*a, **kw):
62     ... func(a, kw, args, kwargs)
63     ... return wrapper
64     ... return real_decorator
65
66     to be written as
67
68     >>> @unwrap
69     ... def decorator(func, *args, **kwargs):
70     ... @wraps(func)
71     ... def wrapper(*a, **kw):
72     ... func(a, kw, args, kwargs)
73     ... return wrapper
74
75     reducing the nested function depth.
76
77     :param dec_func: decorator function
78     :return: unwrapped decorator function
79     """
80     @wraps(dec_func)
81     def unwrap_dec(*args, **kwargs):
82         return wraps(dec_func)(lambda func: dec_func(func, *args, **kwargs))
83     return unwrap_dec

```

Listing B.1: funcwrap.py

REFERENCES

- [1] AKENINE-MOLLER, T., HAINES, E., AND HOFFMAN, N. *Real-Time Rendering*, 3rd ed. CRC Press, 2008.
- [2] BAKER, H. C., AND HEWITT, C. The Incremental Garbage Collection of Processes. *SIGPLAN Not.* 12, 8 (8 1977), 55–59.
- [3] BLEWETT, R., AND CLYMER, A. *Pro Asynchronous Programming with .NET*. Apress, 2014.
- [4] BOWE, N. Hiding IO latency in generators by async prefetching. <http://niki.codekarma.com/2011/05/hiding-io-latency-in-generators-by-async-prefetching/>, 2011.
- [5] BROWN, S., AND VRANESIC, Z. *Fundamentals of Digital Logic with Verilog Design*, 3rd ed. McGraw-Hill, New York City, 2014.
- [6] ECMA INTERNATIONAL. ECMAScript 2017 Language Specification (ECMA-262, 8th edition, June 2017). <https://www.ecma-international.org/ecma-262/8.0/#sec-async-function-definitions>, 2017.
- [7] HALSTEAD, R. H. MULTILISP: A Language for Concurrent Symbolic Computation. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 7, 4 (October 1985), 501–538.
- [8] HOFFMANN, T. PyPI — pythonflow. <https://pypi.org/project/pythonflow/>, 2018.
- [9] JUSTHEURISTIC. PyPI — prefetch-generator. <https://pypi.org/project/prefetch-generator/>, 2016.
- [10] KAPADIA, R., AND BRINZA, B. Getting to the Content You Want, Faster in IE11. <https://blogs.msdn.microsoft.com/ie/2013/12/04/getting-to-the-content-you-want-faster-in-ie11/>, 2013.
- [11] KOSTYUKOV, V. scale.bythebay.io: Vladimir Kostyukov, Futures: Twitter vs Scala. <https://www.youtube.com/watch?v=jiYe-LdPrS0>, 2017.
- [12] KRISHNAMURTHI, S. *Programming Languages: Application and Interpretation*, 2nd ed. Brown University, Providence, 2017.
- [13] LOTT, S. F. *Functional Python Programming: Discover the power of functional programming, generator functions, lazy evaluation, the build-in itertools library, and monads*, 1st ed. Packt Publishing, Birmingham, 2015.
- [14] MIDKIFF, S. P. *Automatic Parallelization: An Overview of Fundamental Compiler Techniques*. Morgan & Claypool Publishers, 2012.
- [15] OHNO, T. *Toyota Production System: Beyond Large-Scale Production*. CRC Press, New York City, 1988, ch. Starting from Need.

- [16] OPENMP ARCHITECTURE REVIEW BOARD. *OpenMP Application Programming Interface Version 4.5*. OpenMP, November 2015.
- [17] PATTERSON, D. A., AND HENNESSY, J. L. *Computer Organization and Design: The Hardware/Software Interface*. Elsevier, Oxford, 2014.
- [18] PRASAD, K., PATIL, A., AND MILLER, H. *Programming Models for Distributed Computing*. <http://dist-prog-book.com/chapter/2/futures.html>, ch. Futures and Promises.
- [19] PYTHON SOFTWARE FOUNDATION. Initialization, Finalization, and Threads. <https://docs.python.org/3/c-api/init.html>, 2018.
- [20] PYTHON SOFTWARE FOUNDATION. PyPI — the Python Package Index. <https://pypi.org>, 2018.
- [21] PYTHON SOFTWARE FOUNDATION. The Python Standard Library: 17.1. threading — Thread-based parallelism. <https://docs.python.org/3/library/threading.html>, 2018.
- [22] PYTHON SOFTWARE FOUNDATION. The Python Standard Library: 17.2. multiprocessing — Process-based parallelism. <https://docs.python.org/3/library/multiprocessing.html>, 2018.
- [23] PYTHON SOFTWARE FOUNDATION. The Python Standard Library: 17.4. concurrent.futures — Launching parallel tasks. <https://docs.python.org/3/library/concurrent.futures.html>, 2018.
- [24] PYTHON SOFTWARE FOUNDATION. The Python Standard Library: 19.5. asyncio — Asynchronous I/O, event loop, coroutines and tasks. <https://docs.python.org/3/library/asyncio.html>, 2018.
- [25] PYTHON SOFTWARE FOUNDATION. The Python Standard Library: 2. Built-in Functions. <https://docs.python.org/3/library/functions.html>, 2018.
- [26] PYTHON SOFTWARE FOUNDATION. The Python Standard Library: 8. Compound statements. https://docs.python.org/3/reference/compound_stmts.html, 2018.
- [27] THE LINUX INFORMATION PROJECT. Pipes: A Brief Introduction. <http://www.linfo.org/pipe.html>, 2006.
- [28] THE MATHWORKS INC. Parallel for loop - MATLAB parfor. <https://www.mathworks.com/help/matlab/ref/parfor.html>, 2018.