

Lazy on Principle

Anonymous Author(s)

Abstract

This is the abstract.

ACM Reference Format:

Anonymous Author(s). 2024. Lazy on Principle. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

When sequences of data become too large to fit into memory at once, programs need to process them lazily. From the humble iterator to asynchronous APIs for streams and sinks with error handling and buffering, every language needs libraries for working with lazy sequences.

For such a fundamental, conceptually simple, and language-agnostic problem, one might expect a principled, unified solution that programming language designers and library authors can turn to and implement in their language of choice.

But the opposite is the case. Learning a new programming language implies learning yet another, slightly (or not so slightly) different set of APIs for working with sequences. Even within a single language, there are often competing libraries — section 8 lists some thirty popular Javascript libraries alone.

Starting from “Which abstraction is the best?”, we quickly moved to “Is there a best abstraction?”, and then to the more constructive “What would make an abstraction the best?”. In this paper, we present our answers to these questions. In a nutshell:

1. Abstractions for working with lazy sequences in the wild are ad-hoc designs.
2. We propose a principled way of evaluating them.
3. No prior abstractions satisfy all evaluation criteria.
4. We develop abstractions that do.
5. ~~Everybody everywhere should use our abstractions without further reflection.~~

Note that we will focus on strictly evaluated languages. This makes explicit the design elements that enable laziness.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2024 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

We further restrict our focus to the two simplest ways of interacting with a (possibly infinite) sequence: *consuming* a sequence item by item, or *producing* a sequence item by item. Both modes of interaction are of great practical interest, they correspond, for example, to reading and writing bytes over a network. We do not consider more complex settings such as random access, or mixing reading and writing.

1.1 Evaluating Sequence APIs

Equipped with a vague notion of wanting to “lazily consume or produce sequences”, how can we do better than simply trying to find a design that satisfies all use-cases we can come up with? In mathematics, one would define a set of criteria that a solution should satisfy, in a way that makes no assumptions about any possible solutions themselves.

For example, a mathematician might want to work with numbers “with no gaps in-between” (i.e., the real numbers). They might formalize this intuitive notion as a minimal, infinite, complete ordered field. Any candidate construction (say, the Dedekind cuts of rational numbers), can now be objectively measured against the requirements. As an added bonus, it turns out that all constructions satisfying the abstract requirements are isomorphic. Some constructions might be more convenient than others in certain settings, but ultimately, they are all interchangeable.

This approach of construction-independent axiomization is the only way we can conceive to bring clarity to the proliferation of competing library designs.

Sadly, we could not find an airtight mathematical formalism to capture our problem space. The criteria we now present leave gaps that must be filled by argumentation rather than proof, the API design still remains part art as much as science. This makes the following paragraphs the weakest link of this paper. We nevertheless think that both our approach and the designs it yields are novel — and useful.

The criteria by which we shall evaluate lazy sequence abstractions are minimality, symmetry, and expressivity.

Minimality asks that no aspect of the API design can be expressed through other aspects of the design. Removing any feature impacts what can be expressed.

Symmetry asks that reading and writing data should be dual. The two intuitive notions of producing and consuming a sequence item by item are fully symmetric and sit on the same level of abstraction. Any API design that introduces an imbalance between the two must either be contaminated with incidental complexity, or it must be missing functionality for one of the two access modes.

Expressivity asks that the API design is powerful enough to get the job done, but also no more powerful than necessary.

This is by far the most vague of our criteria, because we cannot simply equate more expressivity with a better design. We *can*, however, draw on the theory of formal languages to categorize the classes of sequences whose consumption of production can be described by an API. Some of these classes are more natural candidates than others.

Of these criteria, minimality is arguably the least controversial. Symmetry turns out to be the one we generally find the most neglected in the wild, and strict adherence to symmetry shapes the designs we propose to differ significantly from any others we are aware of. Expressivity might have the weakest definition, but turns out to be rather unproblematic: real-world constraints on the APIs lead to a level of expressivity that also has a convincing formal counterpart — the ω -regular languages (see ?? for details) — making us quite confident about the appropriate level of expressivity.

TODO argue for *consistency* as a cure to the painfulness of building async stuff

To obtain a first indicator for an appropriate level of expressivity, we examine the world of non-lazy sequences, i.e., sequences that can be fully represented in memory.

1.2 Case Study: Strict Sequences

Representing sequences in memory can be done in such a natural way that we have never seen any explicit discussion. We shall assume a typical type system with product types (denoted (S, T)), sum types (denoted $S + T$), and homogeneous array types (denoted $[T]$).

Let T be a type, then T is also the type of a sequence of exactly one item of type T . Now, let S and T be types of sequences. Then (S, T) denotes the concatenations of sequences of type S and sequences of type T , $S + T$ denotes the sequences either of type S or T , and $[T]$ denotes the concatenations of arbitrarily (but finitely) many sequences of type T . None of this is particularly surprising, we basically just stated that algebraic data types and array types allow you to lay out data sequentially in memory.

Slightly more interesting is the blatant isomorphism to regular expressions. Each of the “sequence combinators” corresponds to an operator to construct regular expressions; the empty type and the unit type correspond to the neutral elements of the choice and concatenation operator respectively.

This is useful for making our expressivity requirement for lazy sequence APIs more precise: if the natural representation of strict sequences admits exactly the regular languages, then the regular languages are also the natural candidate level of expressivity for lazy APIs.

Unlike strict sequences that have to fit into finite memory, lazy sequences can be of infinite length. The natural generalization of the regular languages are the ω -regular languages. Hence, this is the level of expressivity we want to see in lazy APIs.

The strict case also neatly supports the design goals of minimality and symmetry. Removing any combinator leads

to a strictly less expressive class of languages, and every operator comes both with a way of building up values and with a way of accessing values.

By generalizing the strict case to the lazy case, we can make our requirement of expressivity more precise, leading us to our final set of requirements: We want APIs for lazily producing or consuming sequences an item at a time, such that there is a one-to-one mapping between API instances and ω -regular languages, no aspect of the APIs can be removed without losing this one-to-one mapping, and there is full symmetry between consumption and production of a sequence. Still not entirely formal, but close enough to meaningfully evaluate and design APIs.

1.3 Organization

TODO

2 Related Work

TODO

Iterees: https://link.springer.com/chapter/10.1007/978-3-642-29822-6_15 <https://hackage.haskell.org/package/iterIO-0.2.2/docs/Data-IterIO.html>

session types

pull-stream paper

<http://jiangxi.cs.uwm.edu/publication/rebls2020.pdf>

<https://arxiv.org/pdf/1612.06668.pdf> (good related work)

3 Evaluating Abstractions

We now start by introducing a notation for API designs. We then express several APIs we encountered in the wild in our notation, in order to build intuition, demonstrate sufficient notational expressivity, and to highlight typical violations of our design goals in popular APIs. We do not aim for an exhaustive survey of APIs, we merely select some examples to illustrate our points.

In the following, we use uppercase letters as type variables. (S, T) denotes the product type of types S and T (intuitively, the cartesian product of S and T), and $()$ denotes the unit type (intuitively, the type with only a single value). $S|T$ denotes the sum type of S and T (intuitively, the disjoint union of S and T), and $!$ denotes the empty type (intuitively, the type that admits no value). Finally, we write $S \rightarrow T$ for the type of (pure) functions with an argument of type S and a return value of type T . Note we take a purely functional approach here: a function does not mutate its argument, it simply produces a new value.

We specify an API as a list of named types (typically functions). Each API can quantify type variables that can be used in its *members*¹. As an example, consider the following API:

API *Iterator* $\langle P, I \rangle$

next: $P \rightarrow (I, P) \mid ()$

¹More formally, this is a notation for ad-hoc polymorphism like Haskell’s type classes, Java’s interfaces, or Rust’s traits.

This pseudo-type fragment states that in order to obtain a concrete `Iterator`, one needs two types: a type `P` (**P**roducer) and a type `I` (**I**tem). These types have to be related through existence of a function `next`, which maps a producer to either an item and a new producer, or to a value that signifies that no further items can be produced.

To consume this iterator, one would repeatedly call `next` on the producer returned from the prior call of `next`, until no further Items are produced.

A concrete example of an iterator are the homogenous arrays of `Is` as producers of `Is`; `next` returns `()` for the empty array, otherwise it returns the first item in the array and the array obtained by removing the first item.

This API is completely stateless, we never mutate any `P`. In an imperative programming language, one would typically use a function that takes a reference to a `P` and returns either an `I` or `()`, and then makes all implementors pinky-swear to not invoke the function with a `P` that has returned `()` previously.

We prefer the purely functional notation, because it can express the pinky-swearing API contract on the type level. But all our designs can easily be translated into an imperative, stateful setting. The other way around, by converting stateful references into input values and output values, we can represent APIs from imperative languages in our notation. For example, this `Iterator` API captures the semantics of commonly used iterator APIs such as those of Python² or Rust³. It handily abstracts over the fact that Rust has actual sum types, whereas Python signals the end of iteration with an exception.

Without devoting too much space to it, we want to point out that even for such a simple special case of sequence processing as iteration — i.e., synchronous production of finite sequences without buffering or error handling — there are several, non-isomorphic approaches in the wild. Both Java⁴ and Javascript⁵ use APIs with slightly different properties; they differ in how early the code interacting with the iterator knows that no further items will be produced.

Given the symmetry between producing and consuming data, the virtual non-existence of (synchronous, non-error-handling) APIs for consuming data step by step is jarring. Java has an `Appendable` interface⁶, but it is specialized to characters, and has a much less prominent role than the `Iterator` interface. Clojure complects both production and consumption into a single `Sequence` interface⁷. Several languages do give the opposite of iterators a prominent role: for loops that consume an iterator item by item. Baking

²<https://wiki.python.org/moin/Iterator>

³<https://doc.rust-lang.org/std/iter/trait.Iterator.html>

⁴<https://docs.oracle.com/javase/8/docs/api/java/util/Iterator.html>

⁵<https://tc39.es/ecma262/multipage/control-abstraction-objects.html#sec-iterator-interface>

⁶<https://docs.oracle.com/javase/8/docs/api/java/lang/Appendable.html>

⁷https://clojure.org/reference/sequences#_the_seq_interface

sequence production via first-class values into a language while turning sequence consumption into a purely syntactic component with no run-time presence massively breaks symmetry from a high level design perspective.

Moving beyond synchronous iterators, *asynchronous* producers — often called *streams* — typically come with buffering and error handling. We shall abstract over⁸ both buffering and asynchronicity, as they do not influence matters of minimality, symmetry, or expressiveness regarding formal languages. We believe that designing a solid API irrespective of buffering and asynchronicity and then adding them in a way that is idiomatic for the programming language in question leads to clearer designs. In particular, we see no reason why synchronous and asynchronous sequence APIs in the same language should not be completely analogous.

An example where the asynchronous producer API is equivalent to the synchronous iterator API is — as of writing — Rust⁹. This is rather atypical, because asynchronous APIs, which are often motivated by networking, usually emphasize error handling. In the iterator API, one can use a sum type of actual items and an error type as the type `I`, but on the type level, it remains possible to continue iterating after an error. The API is accurate for recoverable errors only.

A different example is Swift¹⁰, which superficially appears to offer an equivalent API, but the language-level feature of throwing exceptions allows the ability to express irrecoverable errors¹¹. In our notation, the Swift API is:

```
API FallibleIterator <P, I, E>
  next: P -> (I, P) | () | E
```

This particular API design violates minimality: removing the option of returning the unit type would leave the degree of expressivity unchanged, since one could always instantiate `E` as a sum type of the unit type and an actual error type. It might be tempting to argue that the Swift API communicates intent more clearly, and allows for nicer library functions for working with streams. But these conveniences and specializations could just as well be implemented as special cases of the more general, underlying pattern. Rust nicely demonstrates this by offering a host of functions¹² for working with streams whose items are a sum type of actual items and error

⁸The seasoned Haskell-afficionado will immediately see that parameterizing our presentation over a monad for effectful computation [Wad95] would restore rigorous reasoning to our handwavy act of abstraction. We posit that to the readers who are already familiar with effect-management through monads, filling in the blanks is an easy exercise. To those readers who are not familiar with this technique — i.e., the vast majority of practitioners we would like to reach — obscuring our presentation behind higher-kinded type constructors poses an unnecessary barrier to access.

⁹<https://docs.rs/futures/0.3.30/futures/stream/trait.Stream.html>

¹⁰<https://developer.apple.com/documentation/swift/asyncsequence>

¹¹<https://developer.apple.com/documentation/swift/asyncthrowingstream>

¹²<https://docs.rs/futures/latest/futures/stream/trait.TryStreamExt.html>

values. Offering a special case as the most fundamental API, like Swift does, unnecessarily reduces expressivity.

Another school of APIs defines asynchronous producers in terms of all the ways in which one would commonly interact with them: mapping, reducing, piping into a consumer, etc. Examples include Java¹³ or Dart¹⁴. Such designs are not concerned with minimality at all. Our preferred approach is to find the minimal interface that that allows expressing all these higher-level functions on top.

The story for asynchronous *consumers* is usually highly asymmetrical again. Swift, for example, has asynchronous for loops to consume streams. Rust offers a *sink* abstraction¹⁵ as the opposite of a stream, but the designs are highly asymmetric — the sink abstraction requires four functions to the stream abstraction's one function.

4 A Principled Design

From this backdrop of rather inconsistent designs, we now present our alternatives. We derive — in tandem — APIs for producing and consuming sequences, guided by minimality, symmetry, and expressivity (section 4.1). Next, we argue that the designs are indeed sufficiently symmetric and of appropriate expressiveness (section 4.2) — if our arguments are not fully formal, they are at least *formalizable*. Finally, we sketch some consequences the designs could have for language-level features such as generators or for loops (section 5.2).

4.1 Deriving Our Design

To derive a principled design step by step, we start with simplest producer API: a producer that emits an infinite stream of items of the same type.

```
API InfiniteProducer <P, I>
  produce: P -> (I, P)
```

An iterator, in contrast, expresses a *finite* stream of items, by making the result of a sum type with a unit type option to signal termination. We can easily abstract over both through a simple realization: we can rewrite the produce function of the InfiniteProducer as a sum with the empty type, without changing the semantics at all (it is impossible to provide an instance of the empty type, the empty type is the neutral element of type sums):

```
API AlsoAnInfiniteProducer <P, I>
  produce: P -> (I, P) | !
```

Now, the infinite producer and the finite producer have the exact same form, and we can introduce a type parameter for the summand to express either:

```
API Producer <P, I, F>
```

```
  produce: P -> (I, P) | F
```

Setting the type parameter *F* (for *final item*) to *!* or *()* yields the infinite streams and the finite streams over a single type of items, respectively.

Another natural choice for *F* are irrecoverable error types; most APIs with this design denote the type parameter as a type of errors explicitly. This denotation obscures how the same abstraction can also represent iterators or infinite streams, however.

Moreover, it obscures that *F* might be another producer itself, with which to continue production. Through this use of the API, we can effectively concatenate any producer after any finite producer. This usage is the cornerstone of achieving the expressivity of the ω -regular languages, and one we have not encountered in the wild at all.

To give a tangible example of how this degree of expressivity can be useful, consider a networking protocol that proceeds in stages: first a handshake for connection establishment, followed by an exchange of key-value pairs that signify the capabilities of an endpoint, followed by the application-level message exchange. With an API parameterized over arbitrary final values, you can implement each stage in a type-safe way, and then concatenate the stages both in execution and on the type-level. Traditional APIs force programmers to either lump the different kinds of messages (handshake, key-value pairs, application-level) into a single sum type, or to forego helpful typing altogether and operate on the level of bytes.

A symmetric consumer API should be one that can be given either an item of type *I* — returning a new consumer value to continue the process — or a final item of some type *F* — without returning a consumer to continue with. Ideally, we should be able to mechanically derive this API as a dual of the producer API. A tempting option is to “flip all arrows” and simply swap argument and return type of the produce function:

```
API Recudorp <P, I, F>
  ecudorp: ((I, P) | F) -> P
```

We can clean this up by splitting the function of a sum type argument into two independent functions (both versions are equivalent), and giving more conventional names:

```
API NotQuiteConsumer <C, I, F>
  consume: (I, C) -> C
  create: F -> C
```

Unfortunately, this does not give the kind of API we were hoping for. The consume function is appropriate, but the second function is not closing a consumer, but creates a consumer. A straightforward dual construction gives too strong of a reversal to yield an API suitable for practical use.

Instead of a fully dual construction, we instead derive a consumer API in steps analogous to those for deriving the

¹³<https://docs.oracle.com/javase/8/docs/api/java/util/stream/Stream.html>

¹⁴<https://api.dart.dev/stable/3.3.0/dart-async/Stream-class.html>

¹⁵<https://docs.rs/futures/0.3.30/futures/sink/trait.Sink.html>

Producer API. We start again with the consumers of an infinite sequence and the consumer of a finite sequence:

```
API InfiniteConsumer <C, I>
  consume: (I, C) -> C
```

```
API FiniteConsumer <C, I>
  consume: (I, C) -> C
  close: C -> C
```

We can again introduce a type parameter for the final sequence item to unify the APIs; observe how using ! or () for the parameter F in the following API yields results isomorphic to the InfiniteConsumer and FiniteConsumer APIs respectively:

```
API Consumer <C, I, F>
  consume: (I, C) -> C
  close: (F, C) -> ()
```

This API is fully symmetric to the Producer: the consumer can consume exactly those sequences that a producer can produce. It is rather unusual in that we have never seen an API whose close function takes an argument in the wild.

Another unusual aspect is the inability of a consumer to report errors to its calling code. This is severe enough of a departure from typical APIs that we discuss this in more detail in section 4.3. But first, we will argue that these Producer and Consumer APIs indeed fulfil our initial design requirements.

4.2 Evaluating Our Design

Convincing arguments for the symmetry of the APIs are not immediately apparent; while we derived both APIs through analogous steps, they are not immediately dual in any obvious sense, and they even have a different number of functions. The closest we could come to a formal argument is to show that they compose in a satisfying way.

Composing a producer with a consumer amounts to piping the data that the producer produces into the consumer:

Require: P, C, I, F are types such that $\text{Producer} \langle P, I, F \rangle$ and $\text{Consumer} \langle C, I, F \rangle$

```
procedure PIPE( $p : P, c : C$ ): ()
  loop
     $x \leftarrow \text{PRODUCE}(p)$ 
    if  $x$  is of type  $F$  then
      CLOSE( $x, c$ )
      return ()
    else
       $c \leftarrow \text{CONSUME}(x.0)$ 
       $p \leftarrow x.1$ 
    end if
  end loop
end procedure
```

Notice that the pipe function returns the unit type, it leaks no information to the outside. More traditional APIs have no way to funnel the final (usually irrecoverable error) value of the producer to the consumer, so the pipe function has to return that value to the outside. Similarly, APIs that allow consumers to emit errors also have to forward such errors from the pipe function.

On a purely abstract level, composing to the unit type evokes the concept of an element and its inverse composing to an identity element. This seems as strong a formal notion of symmetry we can hope for, without *actually* formalizing things.

More concretely, this means that a producer and a consumer together can fully describe a (sub-)program that processes a data stream. While traditional APIs also allow building up programs by combining and modifying producers and consumers, there always has to be ad-hoc glue code for handling errors, and for moving from one stage of processing to the next. Our APIs capture a fuller set of tasks of such programs. In principle, with an expressive library for constructing producers and consumers, it should be possible to describe much wider classes of programs by piping a single producer into a single consumer, with no glue code or ad-hoc error handling. In particular, progressing from one processing stage to the next can be done by having the producer emit a new producer as a final value, and the consumer pipe this value into an inner consumer in its close function.

We can make a similar argument for composing the other way around: it should be possible to create a pair of a consumer and a producer such that the producer produces everything that the consumer consumes, in the same order (an in-memory queue). This is, in some vague sense, the neutral element of transformation steps in a pipeline (we return to this concept in section 5.1). Here again, we see the benefits of the close function taking an argument: we can map this argument directly to the final value to be emitted by produce.

Libraries without this feature usually only offer queues parameterized over a single item type. One common task where this becomes unnecessarily restrictive include adding intermediate queues for buffering that can also delay propagating an error to the queue's consumer until all items that were emitted before the error have been consumed. Another usecase is testing: if you want to test a component by replacing its input producer with a stub, you cannot use the normal in-memory queue as a replacement, because it cannot emit final values (i.e., errors, in most libraries) on demand.

Having argued that our design is indeed symmetric in a meaningful way, we turn to the question of expressivity. Our core argument rests on the observation that each Producer (or Consumer) defines a formal language over an alphabet of atomic types. More precisely, a $\text{Producer} \langle P, I, F \rangle$ emits an arbitrary number of repetitions of values of type I , followed

by a single value of type F . In more traditional notation of a language as a set of strings, it denotes the set $\{I\}^* \circ \{F\}$.

Given this mapping from sequence APIs to languages, which class of languages do our APIs describe? We claim they — in concert with sums, products, and functions — describe the union of the regular and the ω -regular languages.

The ω -regular languages over Σ are the sets of infinite strings over Σ that are either a concatenation of infinitely many words from the same regular language¹⁶ over Σ (infinite iteration), or the concatenation of a regular language and an ω -regular language over Σ , or a choice between finitely many ω -regular languages over Σ .

As already argued in section 1.2, sum types and product types correspond to choice and concatenation of regular expressions respectively. Unlike the strict case, we cannot rely on homogeneous arrays to act as the counterpart to the Kleene operator, but this appears to be where the Producer API comes in (everything applies analogously for the Consumer API): a $\text{Producer}\langle P, I, F \rangle$ can produce an arbitrary number of repetitions of I s, followed by a single F . In particular, $\text{Producer}\langle P, I, () \rangle$ corresponds to the Kleene operator, and $\text{Producer}\langle P, I, ! \rangle$ corresponds to infinite iteration.

Unfortunately, this simple perspective is not fully accurate. Product types as concatenation are too powerful for us: consider a product (P_1, P_2) , where P_1 is a $\text{Producer}\langle P, I, ! \rangle$. The corresponding language would be a concatenation with an ω -language on the left, but this is explicitly ruled out by the definition of ω -regular languages. Another facet of the same problem is that the type (S, T) is not one that describes *first* emitting an S and *then* a T , as it presents both simultaneously.

To solve this, we can restrict the set of well-formed sequence types we consider to pairs $(S, () \rightarrow T)$ for (sums of) non-repeated types S and arbitrary types T , and $\text{Producer}\langle P, S, T \rangle$ for repeated types S . This removes the ability to express concatenations with an ω -language as the left operand, and introduces the required indirection to express “first S , then T ” (remember that we assume our functions to abstract over effects, so there might well be asynchronicity involved in obtaining the T after processing the S).

We shall not dwell on this subtlety any further, because it does not affect our two main points: our API is expressive enough to describe regular (ω -)regular languages, whereas a more traditional API *without* a dedicated type for the final item is *not* expressive enough, resulting in an unjustified reduction in expressive power compared to representing strict sequences in memory.

¹⁶We assume familiarity with regular languages, for an introduction see [HU69], for example. Or do the sensible thing of searching for “regular language” on Wikipedia.

4.3 Communication Flow

As already mentioned, our proposed Consumer API has no way of emitting errors to its calling code, a design decision that appears to make it unsuitable for network programming, for example. More generally speaking, code interacting with a *consumer* can pass information *to* the consumer but cannot obtain any information *from* the consumer; and conversely, code interacting with a *producer* can obtain information *from* the producer but cannot pass any information *to* the producer.

This rigorous restriction on communication flows evokes design choices such as the unidirectional communication primitives of security-focussed micro-kernels like seL4 [MMB⁺13], so there clearly is a place for such constrained APIs. But the consumer API does not seem appropriate as a general-purpose API.

Trying to add the missing communication flows raises some interesting questions. Should consumer return the next consumer *and* another piece of information, or the next consumer *or* another piece of information? What about `close` — should it be able to return extra information, or not? How should symmetry be preserved — does the Producer API require a stop function that lets the surrounding code communicate that (and why) produce will no longer be called? Should produce itself take a piece of information as input?

Our fairly principled approach of aiming for minimal, symmetric, regular-language APIs provides no guidance here, as these communication flows are not part of the underlying problems. Any choices we need to make are essentially arbitrary.

We see two ways out of this problem. The simplest solution is acceptance. When a programmer wishes to write data to a network through a consumer interface, they need a corresponding producer to emit any feedback such as connection failures. Considering that typical networking APIs use the same error type for reading and writing data, this doesn't seem too far-fetched. Then again, the difficulties in migrating from more typical APIs to this style of error handling are hard to estimate. An argumentative essay like this one cannot conclusively establish a result, we merely want to raise that accepting a consumer API without error reporting might be more feasible than it appears at first glance.

The other solution is to consider fallibility as an *effect*. Just like the functions we use in our APIs might be asynchronous, the might also be fallible. Different programming languages represent the different ways: some could use exceptions, others could consistently use a `Result` type (a sum type of either the actual value of interest or an error value) — the latter is a simple and classic example of monadic effect handling. We can keep using the same notation as before, but consider every function as possibly fallible.

To make explicit the transformation from our APIs to fallible versions, here is what the results look like errors of

type E , when neither relying on exceptions nor on explicitly abstracting over effects via monads:

```
API FallibleConsumer <C, I, F, E>
  consume: (I, C) -> C | E
  close: (F, C) -> () | E
```

```
API FallibleProducer <P, I, F, E>
  produce: P -> (I, P) | F | E
```

But interesting to consider that the `FallibleProducer` API appears asymmetric to the `FallibleConsumer` API without the context of errors-as-effects, as it does not mirror the communication flow of the consumer by having functions that take E s as arguments. The return type of `produce` also appears to violate minimality, as E and F could be combined into a single type parameter in principle. This invalidates our earlier criticism of the Swift-derived `FallibleIterator` API of section 3: a return type of $I \mid () \mid E$ correctly separates the final item (hardcoded to type $()$ in the Swift case) from the error effect of type E .

It took a few pages, but now we have a more nuanced view by which to understand the API designs encountered in the wild. This deeper understanding can hopefully serve to guide future API designs that otherwise would have to rely on gut-feeling and copying older designs.

We will continue our investigation with the `Producer` and `Consumer` APIs as they are, it is up to the reader to decide whether their functions should be fallible (and/or asynchronous, for that matter), or not.

5 Working With Producers and Consumers

Having settled on designs for `Producer` and `Consumer` APIs, we now turn to the can or should be used in practice. We note a powerful pattern composability in section 5.1, muse about language-level support in section 5.2, before turning to matters of efficiency in section 5.3.

5.1 Conducers

In section 4.2, we briefly considered an in-memory queue: a pair of a consumer and a producer such that the producer emits exactly the item consumed by the consumer. We can consider such a pair as a single value that implements both the `Consumer` and the `Producer` API; we shall call such a value a *naïve conductor* (portmanteau of *consumer* and *producer*). The *naïveté* will become apparent once we go from intuitive notions of composability to actual implementation; for now we ask you to suspend some disbelief and let the concept guide us to the more useful, *actual* conductors.

Naïve conductors make an appealing foundation for constructing and composing producers and consumers. You can use a single naïve conductor definition to both obtain a new producer from a producer or a new consumer from a consumer. Consider the naïve queue conductor: composing a

producer with the naïve conductor yields a new producer that buffers some number of items before emitting them. Composing the naïve conductor with a consumer yields a new consumer that buffers some number of items before consuming them in the inner consumer.

This dual-purpose usage constitutes a tangible advantage of being hellbent on symmetry. As a second example, consider a naïve conductor constructed from some function of type $S \rightarrow T$ that is a consumer for items of type S and a producer for items of type T . This naïve *map* conductor can both adapt the items emitted by a producer, or adapt the items accepted by a consumer.

Naïve conductors need not preserve a one-to-one mapping between consumed items and produced items. The common tasks of encoding and decoding values for transport can be captured elegantly by naïve conductors: a *decoder* consumes items of some type S (often, S would be the type of bytes) and occasionally produces an item of some type T , an *encoder* consumes items of some type T and produces many items of some type S .

Unfortunately, none of this actually works. In order to, for example, compose a naïve conductor in front of a consumer, the *consume* function of the resulting consumer would have to first call the *consume* function of the naïve conductor. Then, it would need to correctly guess how many times to call the naïve conductor's *produce* function, in order to feed the results to the inner consumer. A general-purpose composition routine can neither know how many items the inner consumer expects, nor how many items the naïve conductor can produce at any point in time.

One obvious solution is to explicitly manage metadata about which functions can and should be called at runtime, but this creates computational overhead. Another simple solution is to restrict naïve conductor to producing exactly one item per item they consume, but this severely restricts expressivity.

Toward a zero-overhead, expressive solution, we temporarily abandon the dual-usage intuition behind naïve conductors, and examine consumers and producers separately. We define a *consumer adapter* as a function that maps an arbitrary consumer to another consumer, and a *producer adapter* as a function that maps an arbitrary producer to another producer.

These adapters can implement the same functionality as naïve conductors in a way that actually works. Consider, for example, a consumer adapter for encoding items of type S to many items of type T . The consumer adapter can produce a consumer that consumes an item of type S , computes the encoding, and calls the *consume* function of the inner consumer once for each T of the encoding. The corresponding producer adapter, when asked to produce a value of type T , asks the wrapped producer for value of type S , and computes the encoding. It then returns the first T of the encoding and buffers the remaining encoding, to be admitted on subsequent calls

to *produce*. Only when the buffer has become empty does it request another item from the wrapped producer.

There is a large amount of overlap or symmetry between the encoding consumer adapter and the encoding producer adapter, note how both use the same procedure for the actual encoding, and both need to buffer the result in between subsequent calls to the wrapped consumer producer. We call a pair consumer and producer adapters that implements a naïve conductor an (actual) *conducer*.

While such conductors are an interesting tool to reason about working with lazy sequences, they do not provide an immediate software engineering benefit: the two adapters need to be implemented independently. In the spirit of full symmetry, we now have to duplicate all implementation efforts.

To improve on this, we next take a look at how programming language syntax (or macros) to make it possible to write a single definition that then yields both adapters of a conductor. To do so, we first need to investigate dedicated syntax for producers and consumers separately.

5.2 Syntax Considerations

Many programming languages offer generator syntax for creating iterators, and for loops for consuming iterators. A language designed with our APIs in mind could provide more powerful syntax.

Generators¹⁷ provide dedicated syntax for creating producers, with *yield* emitting repeated items and *return* emitting the final value. As an example, the following pseudo-code emits the numbers from zero to nine and then the final string “hi”. We use atypical choices of keywords (producer instead of generator, produce instead of yield, and produce final instead of return) to be obnoxiously explicit about the intended semantics, and to prepare for a symmetric consumer design:

```
producer
  i = 0
  while i < 10
    produce i
  produce final "hi"
```

We are not aware of any language that provides a symmetric construction for creating consumers. Dreaming up a symmetric design initially seems straightforward enough:

```
consumer
  x = consume
  y = consume
until consume final z
  doSomething(x + y + z)
```

This design does leave open some questions: what if the consume function of the created consumer is called more

often then there are consume keywords in the main consumer body? And should it always be valid to jump to the `until consume final` block, or only at the end of the main consumer body?

Since the basic consumer design allows no communication to the calling code, a simple solution to the problem of too many *consume* calls is to implicitly wrap the main consumer body in a loop. In a setting with fallible consumers, a consumer that wants to limit the number of possible calls to *consume* can simply add an extra consume expression and throw from there:

```
consumer
  x = consume
  y = consume
  _ = consume
  throw "too_much_information"
until consume final z
  doSomething(x + y + z)
```

To allow for control about what to do when *close* is called depending on the current state of the consumer, the naive `until consume final` can be replaced with a mechanism that mimics try-catch blocks:

```
consumer
  consumeblock
    x = consume
  until _
    throw "too_little_information"
  consumeblock
    y = consume
  until z
    doSomething(x + y + z)
  consumeblock
    _ = consume
  until _
    throw "too_much_information"
```

The syntax is deliberately painful: we do not claim that these are the best design choices, we merely want to demonstrate that providing a meaningful and useful consumer syntax is indeed possible. And after extrapolating the logic that leads to our API designs, designing generators into languages without a corresponding consumer equivalent feels questionable.

A particular usecase we want to highlight for explicit (asynchronous) consumer syntax is that of implementing asynchronous parsers. Typically this involves writing a state-machine or otherwise putting a lot of manual work into ensuring a parser that can suspend its execution when reaching the temporary end of input and then resume once more input becomes available. The consumer syntax allows writing asynchronous parsers that look just like synchronous ones.

¹⁷<https://peps.python.org/pep-0255/>

Assuming the questions around dedicated consumer syntax have been solved, the next logical step is to combine the consumer and producer keywords into a more powerful conductor language construct. As an example, we sketch an encoder conductor for 16-bit integers into 8-bit integers:

```
conductor
  consumeblock
    x = consume
    produce x / 256
    produce x % 256
  until f
    produce final f
```

From such a construct, both the consumer adapter and the producer adapter can be generated. For the consumer adapter, the *consume* expressions provide the entry points to the state machine of the *consume* function, and each *produce* expression translates to a *consume* call of the wrapped consumer. For the producer adapter, the *produce* expressions provide the entry points to the state machine of the *produce* function, and each *consume* expression translates to a *produce* call of the wrapped producer.

Finally, we want to draw a parallel to coroutines[MI09], as implemented, for example, in Lua[Ier06]. In (that particular brand of) coroutines, the *yield* expression in the coroutine implementation not only yields a value to the outside world, but it also evaluates to a value that is given as part of the expression that resumes the coroutine. We can see our conductor syntax as a generalization of this pattern. Coroutines tie incoming and outgoing communication to the same points in the coroutine, marked by *yield*. In fact, this is equivalent to naïve conductors restricted to maintaining a one-to-one correspondence between consumption and production. Our syntax allows arbitrarily splitting the communication. Hence, conductors generalize coroutines.

5.3 Buffering and Bulk Processing

We now briefly discuss some questions of efficiency. While consumers and producers make for nice building blocks of programs because they are easy to reason about, it is inefficient in practice to process items one by one.

One problem of processing items one at a time is that performing side effects is often expensive, for example, when system calls are involved. Writing a file byte by byte with individual system calls is orders of magnitude slower than buffering bytes sequentially in memory and writing many bytes with a single system call.

An easy solution is to allow consumers to buffer items internally, leaving them the freedom to arbitrarily delay actual processing indefinitely to optimize for efficiency. When writing to a consumer in order to perform side-effects, the programmer needs a way to force the consumer to stop delaying, *flush* its buffer, and actually trigger the effects:

```
API BufferedConsumer <C, I, F>
  consume: (I, C) -> C
  close: (F, C) -> ()
  flush: C -> C
```

The buffered consumer with a *flush* function is a staple of real-world APIs. The analogous functionality for producers, however, is one we have never encountered. The opposite of *flushing* as much data as possible *out of* a buffer is *slurping* as much data as possible *into* a buffer.

```
API BufferedProducer <P, I, F>
  produce: P -> (I, P) | F
  slurp: P -> P
```

Unlike flushing consumers, slurping producers does not immediately serve to trigger effectful production of items. Still, there are arguments in favor of a *slurp* function on producers that go beyond the consistency gains of maintaining symmetry (although, to be clear, those alone would already suffice in our opinion). Consider a producer that emits items from some effectful source which might stop working at any moment (i.e., a network connection). Slurping allows the programmer to pre-fetch data even though processing the available data might be time-consuming and not yet finished.

System calls are not the only reason for processing data in bulk. Simply copying consecutive bytes in memory from one location to another is significantly more efficient than copying each byte individually. Hence, many programming languages offer APIs for producing or consuming many items at a time by way of *slices* (a pointer paired with the number of items stored consecutively in memory starting at the pointed-to address).

A typical example of such *readers* (producers of many bytes simultaneously) and *writers* (consumers of many bytes simultaneously) are the *Reader*¹⁸ and *Writer*¹⁹ abstractions of the Go language. To translate them into pseudo-types, we write *&r[T]* for a slice of values of type *T* that may be read but not written, and *&w[T]* for a slice of values of type *T* that may be written but not read. The Go APIs then translate to the following:

```
API Reader <R, I, E>
  read: (R, &w[I]) -> (R, Nat) | E
```

```
API Writer <W, I, E>
  write: (W, &r[I]) -> (W, Nat) | E
```

The *read* function *writes* (produces) some number of items into a slice, and returns how many items were written. The *write* function *reads* (consumes) some number of items from a slice, and returns how many items were read. A return value

¹⁸<https://pkg.go.dev/io#Reader>

¹⁹<https://pkg.go.dev/io#Writer>

of zero typically indicates the end of the sequence²⁰. We can easily generalize to arbitrary final values of some type F by requiring the returned number to be non-zero and extending the return sum type by a third²¹ option of type F .

Setting aside the interesting naming choices and the fact that most languages unnecessarily specialize the item type to that of 8-bit integers, these APIs display a perfect symmetry that APIs for operating on individual items usually lack.

It is tempting to think of readers and writers as *generalizations* of producers and consumers respectively, but that viewpoint brings a problematic amount of freedom — which parts should be generalized, and which parts should stay the same? Consider, for example, our restrictions to exclusively reading or writing from slices. This is more restrictive than allowing arbitrary access to the slices, and given the defaults of programming languages (no mainstream languages support write-only pointers), the default choice of many would be unrestricted access to the slices. The Rust community has had to put a lot of energy into dealing with the consequences of such an oversight in its standard library²².

Instead, we propose to think about readers and writers as optimization details: any *read* must be equivalent to a series of zero or more calls to *produce*, and any *write* must be equivalent to a series of zero or more calls to *consume*. This viewpoint precisely defines the semantics of the reader and writer APIs, and cleanly specifies answers to questions that might otherwise be non-obvious: may *read* access the contents of the slice? No. What should *read* or *write* do when given an empty slice? Nothing. Is every (buffered) reader or writer a (buffered) producer or consumer respectively? Absolutely.

This last question is crucial: readers are subtypes of producers, and writers are subtypes of consumers. If you take away only one point from this essay, this is the one. Readers and writers stem from file system abstractions, the duality of reading and writing to or from a file make their symmetry an obvious requirement. Streams and sinks trace back to iterators, which arose from traversal of (polymorphic) data structures, hence making the genericity of items an obvious requirement. If programming languages had routinely linked the two abstractions by a subtyping relation, we could have had fully symmetric, fully generic, unified APIs for decades. Instead, these abstractions have remained incomplete, and, consequently, interoperate badly.

²⁰In an asynchronous setting, if no data is currently available but there might be more data in the future, the functions should block instead of returning zero. In an asynchronous setting, the functions should be parked to be resumed at a later point.

²¹Or a *second* option, if we consider the error case as an effect.

²²Rust allows for uninitialized memory, but *reading* from uninitialized memory is unsafe. See <https://github.com/rust-lang/rfcs/blob/master/text/2930-read-buf.md> and <https://blog.yoshuawuyts.com/uninit-read-write/> for details on how this affects its reader API.

One problem with the reader and writer APIs is that they do not compose very nicely: in order to move data from a reader to a writer, you need to use *allocate* an array into which to first copy the data via *read*, and from which to then copy the data via *write*. An alternate API choice without this problem is to *expose* slices of *internal* buffers instead of *processing* slices of *external* buffers:

```
API BulkProducer <P, I, F>
  producer_slots: (P) -> &r[I] | F
  process_produced: (P, Nat) -> P
  extends BufferedProducer <P, I, F>
```

```
API BulkConsumer <C, I, F>
  consumer_slots: (C) -> &w[I]
  process_consumed: (P, Nat) -> P
  extends BufferedConsumer <C, I, F>
  # To close, use the BufferedConsumer
  # close function
```

The *consumer_slots* function provides a slice into an inner buffer of a BulkConsumer, into which the calling code can write. To trigger actual processing of the written items, the *process_consumed* function notifies the consumer how many items were written and tasks it to consume them. The semantics of calling *process_consumed* with some argument n must be those of calling *consume* n times, with the items written to the slice returned by *consumer_slots*. The BulkProducer API works analogously.

Whereas a writer API requires the data to be *consumed* to be in an array, the bulk consumer is required to organize its *internal buffer* as an array. In practice, things are most efficient if *both* sides of the exchange store data consecutively in memory, so we don't expect this shift in responsibility to make a difference to anyone who uses bulk processing for efficiency reasons anyways.

Our APIs are more low-level than the traditional reader and writer APIs. The *read* and *write* functions — we propose to call them *bulk_produce* and *bulk_consume* — can easily be implemented as helper functions that take a slice and copy its contents into the slots exposed by the bulk API.

Given such *bulk_produce* and *bulk_consume* functions, there are now two semantically equivalent ways of piping a bulk producer into a bulk consumer: *pipe_bulk_consume* uses the *producer_slots* of the producer as the slice argument to *bulk_consume* on the consumer, and *pipe_bulk_produce* uses the *consumer_slots* of the consumer as the slice argument to *bulk_produce* on the producer. Neither of these requires allocation of an external buffer to facilitate the communication.

A final, interesting observation on this topic concerns memory safety. In a language with a concept of uninitialized memory that is acceptable to write to but not to read from, a bulk consumer is free to expose a (write-only) slice of

uninitialized memory in its *consumer_slots* function. Whenever *process_consumed* is called thereafter, the consumer can assume that the memory has been initialized. If the calling code is faulty, this can lead to undefined behavior, making the *process_consumed* function *unsafe* in the Rust sense of the word, i.e., it can trigger undefined behavior when its contract is not fulfilled. There is no such problem with the bulk producer API. Thankfully, the *bulk_consume* helper function fully insulates from this source of errors.

 piping and composing without allocations. bulk conducs (and syntax)

 mention vectored io?

6 Onward!

Combinators.

 Conducer syntax: formal transformation into state machines or calculus terms.

 bulk conducs (and syntax)

 Lengths and limits.

 Asynchrony.

 Trees. DAGs? Digraphs?

 Coroutines/generators generalized. Expressiveness?

 Turing machines. Seeking, and “everything is a file”

 Elastic Turing machines.

 Outside movement.

 from sequences to structured data?

 from sequences to sets (should be easier)?

 proper duality?

 Specialization, but from where? Turing machine? Super-coroutine?

 Optimization and Formalization of co-iterators, can research carry over? Reference Compositional Reasoning for Side-effectful Iterators and Iterator Adapters <https://arxiv.org/pdf/2210.09857.pdf>

<http://www.litech.org/~jed/papers/jmatch-popl2006.pdf> Interruptible Iterators

<https://ora.ox.ac.uk/objects/uuid:87a8d92c-017c-4b3f-89a1-ac3056ad51d7/files/m4762057d0a14ed0f7db1ad98e560fa3f> The Essence of the Iterator Pattern

 Iterators: Signs of Weakness in Object-Oriented Languages <https://dl.acm.org/doi/pdf/10.1145/165507.165514>

 Design and Specification of Iterators Using the Swapping Paradigm https://www.researchgate.net/profile/Bruce-Weide/publication/3187667_Design_and_Specification_of_Iterators_Using_the_Swapping_Paradigm/links/09e415117c907dd6da000007

 Design-and-Specification-of-Iterators-Using-the-Swapping-Paradigm.pdf

 Synthesizing Iterators from Abstraction Functions https://dspace.mit.edu/bitstream/handle/1721.1/87058/Jackson_Synthesizing_20iterators.pdf?sequence=1&isAllowed=y

 Segmented Iterators and Hierarchical Algorithms <https://lafstern.org/matt/segmented.pdf>

 New Iterator Concepts <https://svn.osgeo.org/fdo/branches/3.4/Thirdparty/boost/libs/iterator/doc/new-iter-concepts.pdf>

 frp

 coalgebra, coinduction, etc

 for strict sequence typing: Strongly Typed Heterogeneous Collections <http://softlang.uni-koblenz.de/HList/paper.pdf>

7 Conclusion

TODO

 vision of unified, consistent APIs throughout a full language ecosystem. Getting things right is simpler than getting them wrong.

 On a meta note, this essay constitutes evidence that you can share research results with broad applications to a wide range of programming practitioners, without assuming the kind of deep, intuitive familiarity with the Haskell standard library that requires years of practice to obtain. How did that ever become accepted practice in the first place?

8 Appendix A: Javascript Libraries

This list of javascript libraries for working with lazy sequences is intended to demonstrate that there is a clear need for a solid design that people can fall back to rather than reinventing ad-hoc wheels over and over. We list libraries with at least 200 stars on Github, as of February 2024, found by searching Github for “stream”, “observable”, and “reactive”.

- <https://github.com/staltz/xstream>
- <https://github.com/mafintosh/streamx>
- <https://github.com/getify/monio>
- <https://github.com/getify/asyncquence>
- <https://github.com/cyclejs/cyclejs>
- <https://github.com/winterbe/streamjs>
- <https://github.com/winterbe/sequence>
- <https://github.com/pull-stream/pull-stream>
- <https://github.com/dionyziz/stream.js>
- <https://github.com/caolan/highland>
- <https://github.com/kefirjs/kefir>
- <https://github.com/baconjs/bacon.js>
- <https://github.com/cujojs/most>
- <https://github.com/callbag/callbag>
- <https://github.com/paldepind/flyd>

The following libraries do not explicitly define *streams*, but they do work with *observables*. Observables are an abstraction for values that (discretely) vary over time. For most intents and purposes, this is isomorphic to the notion of a stream.

- <https://github.com/reactivex/rxjs>
- <https://github.com/tc39/proposal-observable>
- <https://github.com/zenparsing/zen-observable>
- <https://github.com/vobyjs/oby>
- <https://github.com/adamhaile/S>
- <https://github.com/luwes/sinuous>
- <https://github.com/mobxjs/mobx>

- <https://github.com/fynyky/reactor.js>
- <https://github.com/ds300/derivablejs>
- <https://github.com/elbywan/hyperactiv>
- <https://github.com/component/reactive>
- <https://github.com/mattbaker/Reactive.js>

These libraries exist in addition to language-level or runtime-level APIs such as the following:

- **Node JS Streams**, and their evolution:
 - `streams0`
 - `streams1`
 - `streams2`
 - `streams3`
- **WHATWG Streams**
- **ECMAScript Iterator**
- **ECMAScript AsyncIterator**

References

- [HU69] John E Hopcroft and Jeffrey D Ullman. *Formal languages and their relation to automata*. Addison-Wesley Longman Publishing Co., Inc., 1969.
- [Ier06] Roberto Ierusalimsky. *Programming in lua*. Roberto Ierusalimsky, 2006.
- [MI09] Ana Lúcia De Moura and Roberto Ierusalimsky. Revisiting coroutines. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 31(2):1–31, 2009.
- [MMB⁺13] Toby Murray, Daniel Matichuk, Matthew Brassil, Peter Gammie, Timothy Bourke, Sean Seefried, Corey Lewis, Xin Gao, and Gerwin Klein. sel4: from general purpose to a proof of information flow enforcement. In *2013 IEEE Symposium on Security and Privacy*, pages 415–429. IEEE, 2013.
- [Wad95] Philip Wadler. Monads for functional programming. In *Advanced Functional Programming: First International Spring School on Advanced Functional Programming Techniques Båstad, Sweden, May 24–30, 1995 Tutorial Text 1*, pages 24–52. Springer, 1995.