

函数式编程中的惰性求值机制

黄梓淳

Jun 11, 2026

在处理超大规模数据或者理论上无限的序列时，我们常常会遇到这样一个问题：如果语言默认把所有表达式都立刻算完再交给后续逻辑，那么当数据量超过物理内存，或者序列根本不存在“终点”时，程序就会直接崩溃或者陷入死循环。惰性求值正是为了解决这一矛盾而出现的求值策略。它把“什么时候需要值”这一决策权从编译期或运行期提前阶段推迟到真正消费数据的时刻，从而让开发者可以用一行简洁的表达式完成原本需要写成复杂循环或迭代器的任务。本文将沿着“概念—理论—实现—实践”的路径，系统梳理惰性求值的原理、语言级支持，以及在工程中的权衡。

1 基础概念：什么是惰性求值

惰性求值 (lazy evaluation) 与及早求值 (eager/strict evaluation) 相对，前者只在表达式结果真正被需要时才执行计算，后者则在绑定发生时立即完成求值。最常见的惰性载体包括 thunk、promise、generator 和 lazy sequence。以 Python 生成器为例，当我们写下 `range(10**12)` 时，解释器并没有在内存里开辟一块能装下一万亿个整数的空间，而是返回一个迭代器对象；只有当 for 循环真正取下一个元素时，生成器才计算并产出下一个值。这种“按需生产”的模式既避免了不必要的存储，也让“无穷”数据结构在语法上成为可能。惰性求值最核心的价值在于三点：其一，避免不必要的计算，从而在性能关键路径上节省 CPU 周期；其二，用有限的代码描述理论上无限或大小未知的数据；其三，把数据生产与数据消费在时间和模块上解耦，让各自的逻辑可以独立演化。

2 理论基石： λ 演算与图归约

λ 演算为惰性求值提供了最底层的数学模型。在 λ 演算中，函数调用对应 β 归约 (β -reduction)。若采用“最左最外” (leftmost outermost) 归约策略，相当于先把函数体整体代入，再逐步把参数归约，这正是正常序 (normal order) 求值；若采用“最左最内” (leftmost innermost) 策略，则先把实参算完再代入，对应应用序 (applicative order)。正常序允许把未被用到的参数整个丢弃，从而天然支持短路；应用序则可能做无用功，但通常更易于在机器上实现。

图归约 (graph reduction) 在树归约的基础上增加了共享节点：如果一个子表达式在多处被引用，归约一次即可，所有引用点同时看到结果。这避免了重复计算，也解释了为什么 Haskell 在默认惰性求值下仍能保持合理的性能。举例来说，对于表达式 `let x = expensive() in (x, x)`，图归约只会调用一次 `expensive`，而树归约则会调用两次。

3 语言层面的实现策略

Haskell 把惰性求值作为语言默认策略。编译器把每个表达式包装成一个 thunk，运行时用一个黑洞 (blackhole) 标记正在求值的 thunk，以检测循环依赖；当 thunk 被强制 (force) 后，其结果会覆盖原 thunk，实现一次求值、多次共享。空间洩漏 (space leak) 是这种策略的副作用：如果一个 thunk 持有对上下文的意外引用，GC 就无法回收这部分内存。严格性分析 (strictness analysis) 可在编译期推断哪些 thunk 必然会被立即使用，从而把它们改写为 eager 版本，缓解空间洩漏。

按需调用 (call-by-need) 与按名调用 (call-by-name) 的主要区别在于是否对 thunk 结果做 memoization。前者把求值结果写回 thunk，后续访问直接返回；后者每次都重新求值。Scala 的 lazy val、F# 的 lazy 关键字，以及 Java 的 Stream，都是显式惰性的例子。Python 的生成器本质上是协程，可以在任意位置挂起并恢复状态，因而也能模拟惰性序列。OCaml 提供了 Lazy 模块，允许程序员在严格求值为主的语言里按需引入惰性。

4 典型抽象与模式

最直观的惰性抽象是“流” (stream)。一个流可以表示为 `Cons(head, () => tail)`，其中 tail 是一个 thunk，只有在被访问时才展开。基于这一结构，我们可以定义自然数流：从 0 开始，每次 tail 返回前值加一的 thunk。斐波那契数列同样可以用两个互递归的 thunk 实现：`fib = Cons(0, () => Cons(1, () => zipWith(+, fibs, tail(fibs))))`。埃拉托斯特尼筛法则把“所有数的序列”作为输入，不断过滤掉当前素数的倍数，产出下一个素数。

在控制结构层面，if、and、or 本质上是惰性函数：只有在条件为真时才对 then 分支求值。利用同样的机制，我们可以自己封装 when、unless、try 等流程抽象，而不需要语言内置特殊形式。惰性 I/O 则把文件句柄的读取动作延迟到消费者真正迭代时才发生，从而让资源生命周期与数据流同步，避免过早打开或迟迟不关闭句柄。Haskell 的 pipes、conduit 以及 Java 的 Reactor 都基于这一思想，把 I/O 操作建模成惰性流。

5 实战案例：从零实现一个迷你惰性列表

我们用 TypeScript 实现一个最小惰性列表，接口设计为 `Cons<T>(head: T, tail: () => LazyList<T>)`。类型定义如下：

```
1 type LazyList<T> = { head: T; tail: () => LazyList<T> } | null;
```

head 函数直接返回节点值：

```
1 function head<T>(list: LazyList<T>): T | undefined {  
  return list ? list.head : undefined;  
3 }
```

tail 函数调用 thunk 获得后续列表：

```
1 function tail<T>(list: LazyList<T>): LazyList<T> {  
  return list ? list.tail() : null;
```

```
3 }
```

take 函数把惰性列表前 n 个元素收集到数组:

```
1 function take<T>(n: number, list: LazyList<T>): T[] {
  const result: T[] = [];
3  while (n-- > 0 && list) {
    result.push(list.head);
5    list = list.tail();
  }
7  return result;
}
```

map 与 filter 都返回新的 thunk 链，而不是立即遍历:

```
function map<T, U>(f: (x: T) => U, list: LazyList<T>): LazyList<U> {
2  return list ? { head: f(list.head), tail: () => map(f, list.tail()) } : null;
}
4
function filter<T>(p: (x: T) => boolean, list: LazyList<T>): LazyList<T> {
6  if (!list) return null;
  if (p(list.head)) return { head: list.head, tail: () => filter(p, list.tail()) };
8  return filter(p, list.tail());
}
```

性能对比实验显示：在只消费前 1000 个元素时，惰性列表的内存占用几乎与 n 无关，而严格列表则随 n 线性增长；但若多次遍历同一惰性列表且未做 memoization，每次都会重复计算。把 tail thunk 改写为带缓存的版本，可把多次遍历的代价降到与严格列表相当。

6 工程权衡与最佳实践

何时选用惰性、何时选用严格，可遵循一条经验法则：若数据量不确定或可能无穷，则倾向惰性；若需要随机访问或必须一次性拿到全部结果，则选用严格。调试惰性代码时，堆栈踪迹经常在 thunk 展开处“断层”，导致难以定位原始表达式。Haskell 提供 BangPatterns 扩展，允许在模式匹配处显式收紧求值；Scala 则可通过 -Xcheckinit 在初始化时就触发异常，从而及早暴露问题。

内存分析需要关注 thunk 数量、共享度以及 GC 压力。若一个列表在构造后只被消费一次，thunk 数量与列表长度呈线性关系；若多处共享同一列表，图归约可把空间复杂度降到 $O(1)$ 。在高并发场景下，惰性求值还可能引入锁竞争，因为多个线程可能同时 force 同一个 thunk。

7 生态与工具链

Haskell 生态中最成熟的惰性流库包括 Streaming 与 Pipes，它们把 I/O 与转换都建模为惰性流，并提供资源安全保证。Scala 的 Cats Effect Stream 与 ZIO ZStream 在 JVM 上实现了类似能力，同时兼容响应式规

范。Java 8 的 `java.util.stream` 虽然默认惰性，但只支持有限的中间操作，且终端操作会立即触发求值。JavaScript 的 `lxJS` 与 Python 的 `itertools` 则把惰性抽象带入动态语言，让开发者可以用函数式风格处理大数据集。性能分析可借助 `perf`、`async-profiler` 或 Python 的 `cProfile`，定位 `thunk` 展开带来的额外开销。

回到导语中的例子：当我们面对一个理论上无限的日志流时，惰性求值让“只打印前 10 条错误日志”变成一行 `logs.filter(isError).take(10)`，既不需要手动管理迭代器，也不会把整个文件读入内存。惰性求值把表达力、模块化与性能三者辩证地统一起来，但也要求开发者理解 `thunk`、共享与空间洩漏等底层机制。

延伸阅读可参考 Okasaki 的《*Purely Functional Data Structures*》，深入探讨持久化数据结构与惰性求值的结合；Peyton Jones 的《*The Implementation of Functional Programming Languages*》则系统讲解了 G-machine 与图归约的实现细节。进一步的方向包括响应式编程中的背压机制、并行惰性求值中的工作窃取，以及为特定领域设计的惰性 DSL。

8 附录

完整示例代码已发布在 GitHub 仓库 `lazy-eval-demo`，使用 TypeScript 4.9 与 Node.js 18 测试通过。术语对照：`thunk`—`thunk`，`call-by-need`—按需调用，`graph reduction`—图归约，`space leak`—空间洩漏。