Open in app

## Lucas Godoy

Follow        11 Followers        About
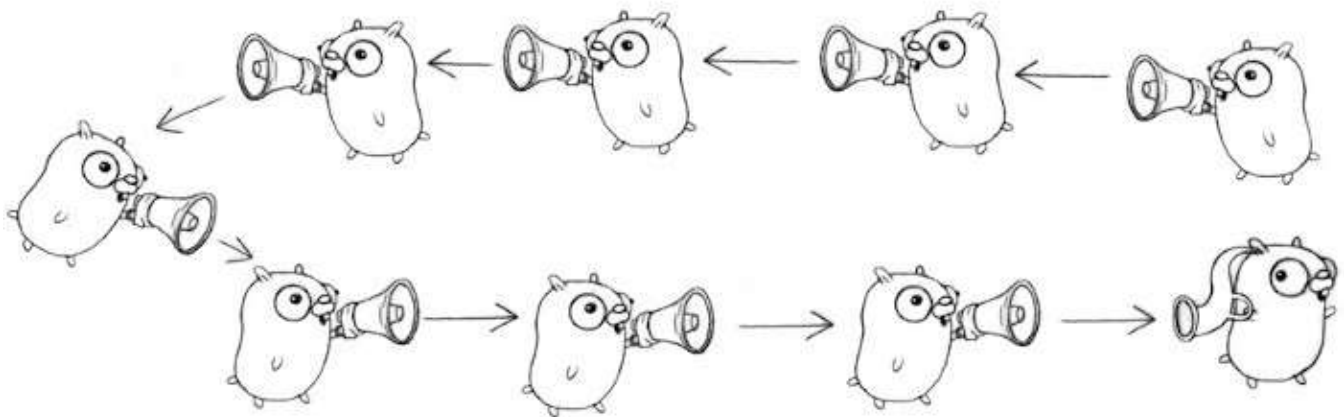
# Go concurrency applied to data pipelines

A different approach to batch processing, and how to potentiate the power of data pipelines throughout the use of the Go concurrency model.

Lucas Godoy  Mar 27 · 6 min read



### Introduction to pipelines

The term **pipeline** applied to the computer science field, is nothing more than a series of stages that take data in, perform some operation on that data, and pass the processed data back out as a result.

Thus, when using this pattern, you can encapsulate the logic of each stage and scale your features quickly by adding/removing/modifying stages, each stage becomes easy to test, nor mentioning the huge benefit of leveraging it by using concurrency, which was the motto for this article.
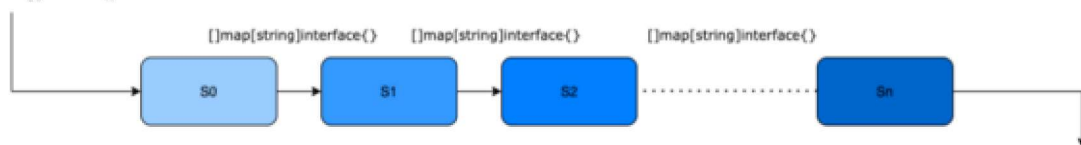
### Previous problem and solution

to integrate retailer's product availability within the main's company app. After
running that integration, the users were able to buy groceries with fewer stockouts
risks.

To accomplish this feature, we built this "availability engine" in GoLang. So, what was
it about? Long story short, it ingested several retailer's CSV files with product
availability information executing a couple of steps to enrich and filter data out based
on certain business logic. At the end of the process, new files were made, with all the
resultant products to be integrated into the company's app for users to buy.

Overall, this data processing pipe looked something like this:



Example of architecture for handling in batches

As the diagram states, the first stage of the pipe takes in a set of CSV rows,
processes them all, and puts the results on a new batch (new slice of maps). And the
very same process repeats it as many times as the number of stages the pipe
actually has. The particularity of this pattern was the fact that none of the next stages

To speed up and optimize the job we start using concurrency at the CSV file level, so we were able to process files concurrently. This approach fits us very well, but there is no silver bullet as we always say...

When working on this project, I gave in the concurrency to one of the more experienced engineers in the team, so I mostly took part in coding each of those business logic steps.

Motivated by the curiosity about different languages leads me to give it a try to the Go concurrency model and patterns to have more tools on my toolbelt. Hence I started reading some posts and books about concurrency in Go. While doing so, I stumbled upon a fantastic pattern, pipelines leveraged by the use of channels, which *if I would have come across it by that time, I would rather suggest implementing it to solve this feature for sure!*

## A better approach for data pipelines: streams of data

In my team's approach, we were using batch processing between stages, which was good enough for us, but there are certainly other options that would have suited better to make it more performant.

Particularly we are speaking of *streaming data* across the different pipeline's stages. This actually means that each stage receives and emits one element at a time instead of waiting for a full batch of results from the previous step to do its thing on them.

In contrast, if we have to compare the memory footprint between batching and streaming, the former is bigger, since each stage has to make a new slice of maps of equal length to store the results of its calculations. On the contrary, the streaming approach will end up receiving and emitting one element at a time so the memory footprint is back down to the size of the pipeline's input.
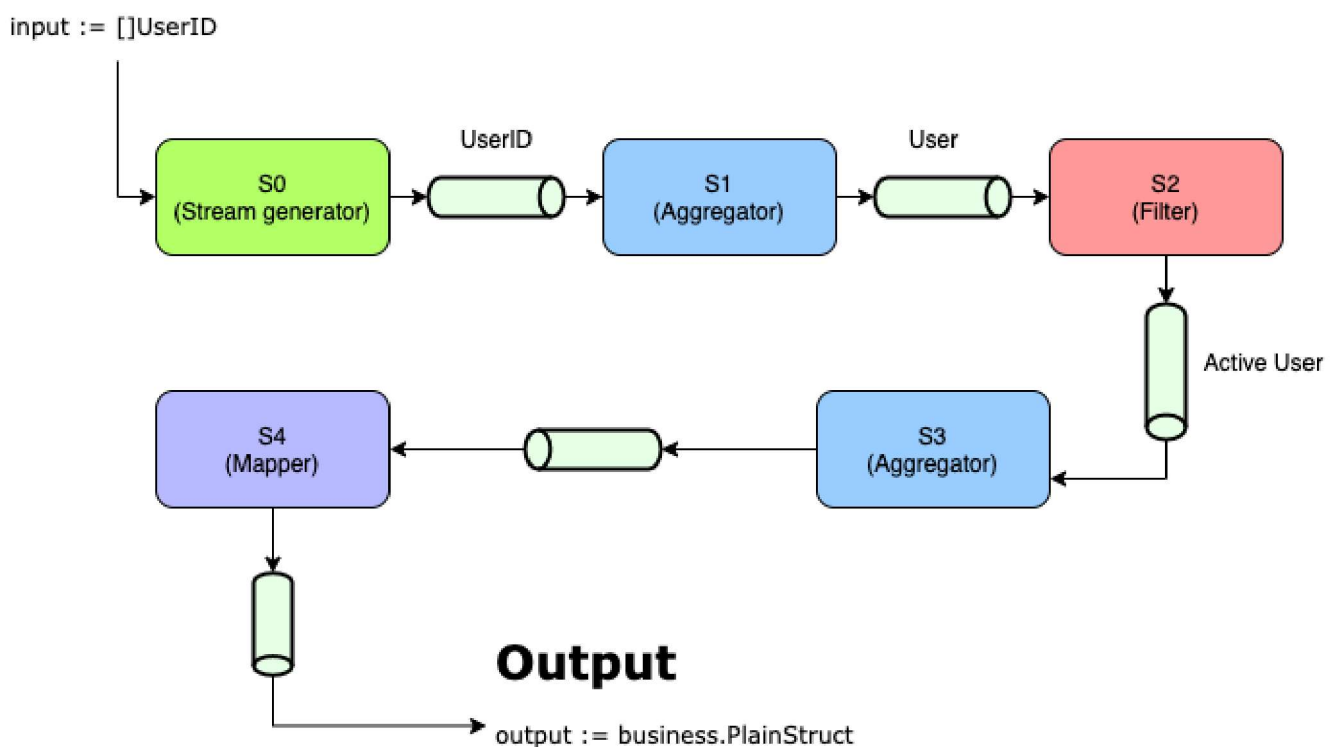
## Implementation example

To demonstrate the pipeline implementation in Go, I had coded a very silly example of this pattern (if you want you can find the source code here). This pipeline is built on top of a chain of stages whose first input would be a series of *UserIDs* (currently just mocked by code). So then, the pipeline stages are the following:

- S2: Filter the inactive users

- S3: Fetch the user's profile and aggregating it into the payload

- S4: Convert the overall aggregation as a plain object to be saved later on somewhere

**Input**

input := []UserID



**Output**

output := business.PlainStruct

Here below, at the *main function,* we can see how the pipeline invocation will look like at line 24.

```go
1    package main
2
3    import (
4            "fmt"
5
6            "github.com/godoylucase/go-data-pipelines-example/business"
7            "github.com/godoylucase/go-data-pipelines-example/pipeline/stages/aggregate"
8            "github.com/godoylucase/go-data-pipelines-example/pipeline/stages/filter"
9            "github.com/godoylucase/go-data-pipelines-example/pipeline/stages/map_from"
10           "github.com/godoylucase/go-data-pipelines-example/pipeline/stages/stream"
11   )
```

Open in app

```go
15   func main() {
16           done := make(chan interface{})
17           defer close(done)
18
19           userIDs := make([]business.UserID, maxUserID)
20           for i := 1; i <= maxUserID; i++ {
21                   userIDs = append(userIDs, business.UserID(i))
22           }
23
24           plainStructs := map_from.UPAggToPlainStruct(done,
25                   aggregate.Profile(done,
26                           filter.InactiveUsers(done,
27                                   aggregate.User(done,
28                                           stream.UserIDs(done, userIDs...)))),
29           )
30
31           for ps := range plainStructs {
32                   fmt.Printf("[result] plain struct for UserID %v is: -> %v \n", ps.UserID, ps)
33           }
34   }
```

Let's start from the innermost function call from our pipeline. The first stage *stream.UserIDs(done, userIDs...)* is the one that will feed the pipeline by streaming the *UserIDs* values. To accomplish this, I used a **generator** pattern, which receives a slice of *UserIDs* (input) and by ranging over it will start pushing each value into a *channel* (output). Therefore, the returned channel will be in turn the input for the next stage.

```go
1    package stream
2
3    import (
4            "fmt"
5
6            "github.com/godoylucase/go-data-pipelines-example/business"
7    )
8
9    func UserIDs(done <-chan interface{}, uids ...business.UserID) <-chan business.UserID {
10           uidsStream := make(chan business.UserID)
11           go func() {
12                   defer close(uidsStream)
```

Open in app                                                                    ●◗

```
15
16                          select {
17                          case <-done:
18                                  return
19                          case uidsStream <- uid:
20                          }
21                  }
22          }()
23          return uidsStream
24  }
```

**pipeline-generator.go** hosted with ♡ by **GitHub**                              view raw

Generator pattern to stream values by using a channel for it

As I mentioned before, conceptually, every stage is fed by input and spits an output out for results. So then, the stages are going to *receive a read-only channel* as input to do the data processing and will *return a read-only channel* on which the resulting outputs will be posted on. Doing so, the next stage on the pipeline will take the previous stage's output as its input for executing its own logic, and so forth till the end of the pipe.

Because of that, the use of channels across the pipe will allow us to safely execute concurrently each pipeline stage because our inputs and outputs are safe in concurrent contexts.

Let's have a look at the following stages on the chain, where based on the streamed data from the first stage (generator) we fetch the actual user data, filter inactive users out, enrich them with its profile, and finally rip some data apart to make a plain object out of the whole aggregation/filtering process.

```
1   package aggregate
2
3   import (
4           "fmt"
5
6           "github.com/godoylucase/go-data-pipelines-example/business"
7   )
8
9   type UserFetcher interface {
10          Get(ID business.UserID) (business.User, error)
11  }
```

Open in app

```go
15          go func() {
16                  defer close(uStream)
17                  for id := range uids {
18                          fmt.Printf("[aggregate] User for UserID %v\n", id)
19
20                          user, err := getUser(id)
21                          if err != nil {
22                                  // TODO address errors in a better way
23                                  fmt.Println("some error ocurred")
24                          }
25
26                          select {
27                          case <-done:
28                                  return
29                          case uStream <- user:
30                          }
31                  }
32          }()
33
34          return uStream
35  }
36
37  // getUser dummy function to simulate some fetching action on some user repository
38  func getUser(ID business.UserID) (business.User, error) {
39          username := fmt.Sprintf("username_%v", ID)
40          user := business.User{
41                  ID:        ID,
42                  Username: username,
43                  Email:    fmt.Sprintf("%v@pipeliner.com"),
44                  IsActive: true,
45          }
46
47          if ID%3 == 0 {
48                  user.IsActive = false
49          }
50
51          return user, nil
52  }
```

pipeline-agg-user.go hosted with ♡ by GitHub          view raw

Fetching users and return them on the channel

```go
1  package filter
2
3  import (
```

Open in app                                                                    ●◗||

```go
6            github.com/godoylucase/go-data-pipelines-example/business
7    )
8
9    func InactiveUsers(done <-chan interface{}, users <-chan business.User) <-chan business.User {
10           activeUsers := make(chan business.User)
11       go func() {
12               defer close(activeUsers)
13               for user := range users {
14                       if !user.IsActive {
15                               fmt.Printf("[filter] Inactive user with UserID %v\n", user.ID)
16                               continue
17                       }
18
19                       select {
20                       case <-done:
21                               return
22                       case activeUsers <- user:
23                       }
24               }
25       }()
26
27       return activeUsers
28   }
```

```go
1    package aggregate
2
3    import (
4            "fmt"
5
6            "github.com/godoylucase/go-data-pipelines-example/business"
7    )
8
9    type ProfileFetcher interface {
10           GetByUserID(uid business.UserID) (business.Profile, error)
11   }
12
13   func Profile(done <-chan interface{}, users <-chan business.User) <-chan business.UserProfileAg
14           upaStream := make(chan business.UserProfileAggregation)
15       go func() {
16               defer close(upaStream)
17               for user := range users {
18                       fmt.Printf("[aggregate] Profile for UserID %v\n", user.ID)
```

Open in app                                                                    ●◗

```go
22                          // TODO address errors in a better way
23                          fmt.Println("some error ocurred")
24                      }
25
26              select {
27              case <-done:
28                      return
29              case upaStream <- business.UserProfileAggregation{
30                      User:    user,
31                      Profile: profile,
32              }:
33              }
34          }
35      }()
36
37      return upaStream
38  }
39
40  // getByUserID dummy function to simulate some fetching action on some profile repository
41  func getByUserID(uids business.UserID) (business.Profile, error) {
42      p := business.Profile{
43          ID:       business.ProfileID(uint(uids) + 100),
44          PhotoURL: fmt.Sprintf("https://some-storage-url/%v-photo", uids),
45      }
46
47      return p, nil
48  }
```

pipeline_profile_agg.go hosted with ♡ by GitHub                          view raw

Aggregating user's profiles to the payload

```go
1   package map_from
2
3   import (
4       "fmt"
5
6       "github.com/godoylucase/go-data-pipelines-example/business"
7   )
8
9   func UPAggToPlainStruct(done <-chan interface{}, upAggregation <-chan business.UserProfileAggre
10      psStream := make(chan business.PlainStruct)
11      go func() {
12          defer close(psStream)
```

```
16                          select {
17                          case <-done:
18                                  return
19                          case psStream <- upa.ToPlainStruct():
20                          }
21                  }
22          }()
23
24          return psStream
25  }
```

**pipeline-map-to-plain.go** hosted with ♡ by **GitHub**      view raw

Converting the payload to a reduced version of it

For sure you have noticed I passed a *done chan interface{}* around the stages. What is this for? It is worth mentioning that goroutines are not garbage collected in runtime, so we have to make sure as programmers to make them all preemptable. By doing that, consequently, we will not leak any goroutine (I will be writing more about this in another post later on) and freeing up memory. Thus, any invocation to the pipeline could be stopped by just closing the *done* channel. This action will lead to the termination of all the spawn children's goroutines and cleaning them up.

All in all, after the latest stage on the pipe, start pushing data out through its output channel another routine (in this case the main one) could start reading out of that channel by iterating over it as stated on main function line 31 and do something with it.

And that is pretty much it! We implemented our first data pipeline, elegant and efficient! (more than batching pipes).

## Final thoughts

In brief, if I any time on I got the chance to work with a similar problem as I had before, I will go for sure for this pattern, which not only is way more performant in terms of memory footprint but also it is faster than doing a batching approach since we can get data being processed concurrently.

Also, there are many other things we can do with pipelines, like rate limiting and fanning it in/out. This topic will be covered later in upcoming posts, where the idea is

This is my first time publishing an article, so any critics are more than welcome, so I can keep learning from all of you guys.

Cheers!

PS: there are a couple of things to be improved in this example that will do make it happen in one upcoming post in which we can talk about error handling for these types of concurrency pattern implementations.

Go      Concurrency      Programming      Data      Pipeline

About   Help   Legal

Get the Medium app