

# XX: A system

Anonymous Author(s)

## Abstract

The promise of serverless currently remains an elusive and extremely attractive paradigm: get what you need as you need it, but only pay for what you use.

For example, building an elastic web server on top of a serverless infrastructure is not something that can be done easily today, but is in principle a good candidate for the serverless setting: its load is unpredictable and bursty, and it does not require managing its data locally.

Many challenges remain in making the ideal of serverless a reality; XX tackles scheduling: in a world where a large and heterogeneous set of jobs is being run on as serverless functions, there has to be a way to differentiate between functions' requirements and urgency.

XX is a distributed scheduler that allows developers to express priorities and enforces them. Developers express priorities to XX via assigning functions to fixed dollar amounts per unit of compute, and cap the overall usage by specifying a monthly budget. Memory costs per unit time used are the same across all priorities, and developers specify a maximum amount of memory for each function. [hmng: Do I need the memory blurb? I added it because I mention memory below ] XX places incoming jobs on machines that have memory available, and implements a machine scheduler that enforces priorities.

[hmng: TODO impl/eval blurb ]

## 1 Introduction

A world where cloud compute is run in the format of serverless jobs is attractive to developers and providers: developers pay only for what they use, while having access to many resources when needed; and cloud providers have control over scheduling and can use that to drive up utilization.

This characteristic of only paying for what you use is especially attractive to developers of applications where the amount of resources that they need varies significantly over time, or is generally small, so that buying their own machines or renting a fixed amount (eg EC2) is untenable.

Some back of the envelope math shows that lambda functions can be cheaper: for a low-traffic website, with approx 50K requests per day, a memory footprint of < 128 MB, and 200ms of execution, that adds up to \$1.58 per month. On the other hand, the cheapest ec2 instance costs just over \$3 per month. Of course, as the number of requests goes up, so does the price for lambda. There are pretty extensive simulations that others have done that show the tradeoff points for different types of workloads. Lambdas can also burst faster:

starting a new lambda execution environment is much faster than starting a new container or ec2 instance, which can take multiple minutes.

There are many reasons that developers choose not to use serverless, despite in theory having workloads that are well-suited for the serverless environment. Some include cold start times, provider lock in, and lack of insight for debugging and telemetry. The issue that this paper addresses is that of the resource requirement interface, and how developers are/are not able to express their needs for resources and latency.

The interface for expressing requirements is already one of the complications of the above modeling: developers are required to choose for each function an amount of memory they will need, which is then tied to a cpu power (an amount of vCPUs). This means that the break-even point and general cost and performance tradeoffs are heavily dependent on the details of the workload, including but not limited to how much memory it will need, how long it will run, and whether it is cpu or i/o bound.

These things are often difficult to know in advance, and more importantly are not correlated with what developers actually care about, which is latency. In fact, measurements have found that in some ways the two are inversely correlated, ie that lambdas that had more memory allocated took longer to start up. This may or may not be outweighed by the speedup associated with more memory, but knowing in advance how much your function will speed up given different amounts of memory can be difficult.

This paper proposes a new interface, that centers priority rather than resource requirements; and a scheduler XX that enforces the priorities laid out. This is getting at the heart of what serverless should be: a pay-per-use model, where resource requirements don't need to be known ahead of time and are expected to change. What is however known, and gets at the core of what developers care about, is priority: some jobs are more important than others, and some clients are willing to pay more money to run their jobs, and fast, than others. Priority is explicitly tied to money in XX: developers express priorities by choosing price classes at which to run their functions.

Centering these classes benefits the developer as well as the provider. The developer is not required to make estimates that may or may not end up being correct, and can just pay for the resources they use. They have access to cheaper options if they are willing to wait, and can skip the line and get high priority access if they pay for it.

Providers, on the other hand, have a concrete metric that can allow them to drive up utilization and profit: they can run the most expensive jobs first, and fill utilization gaps with cheaper jobs.

This also means that there are no clear guarantees when a developer puts a price on a function they want to run. In order to mitigate that somewhat and not go into bidding wars, we propose exposing a fixed set of price classes. This is similar to how AWS has different EC2 instance types, that are directly mapped to prices.

Rather than being a guarantee, the price class is instead a metric to express priority to XX, which it can then use to enforce a favoring of high priority jobs. It is a small change in the API, that still gives the scheduler the information it needs to decide what to schedule when, and aligns the interests of the developer with those of the provider more directly, by communicating on the level of what the provider and developer actually care about: money and latency, as achieved by priority in the system.

[hmng: nowhere in this new intro do I mention memory or challenges; guess its more of a motivation? should that come next?]

A main challenge in designing XX is that of managing memory. Compute can always be divided or preempted, but the buck stops once a machine is out of memory. This problem is made more difficult by the fact that XX does not require memory usage limits to be given by developers, which would otherwise turn memory provisioning into a big packing with overprovisioning problem. XX is thus faced with the challenging proposition of blindly placing jobs not knowing how much memory they will use, but still needing memory utilization to be high.

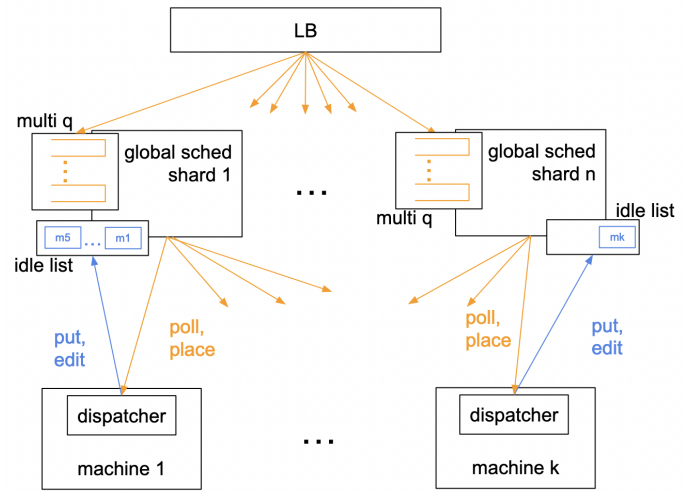
## 2 Design

### 2.1 Interface

Developers using XX write function handlers and define triggers just like they would for any existing serverless offering. In addition, they assign each function to a price class, this is done at function creation. For instance, a simple web server might consist of a home page view that is assigned a higher priority and costs  $2\mu\text{c}$  per cpu second, a user profile page view which is assigned a middle-high priority and cost  $1.5\mu\text{c}$  per cpu second, and finally an image processing job that can be set to a low priority which costs only  $0.5\mu\text{c}$  per cpu second.

Priorities are inherited across call chains: if a high priority job calls a low priority job, that invocation will run with high priority. This is important in order to avoid priority inversion.

Developers pay for memory separately, and by use; the price for memory is the same across all priorities.



**Figure 1:** global scheduler shards queue and place jobs (in orange), on each machine a dispatcher thread keeps track of memory utilization and if it's low writes itself to an idle list (in blue)

To avoid unexpected costs in the case of for example a DOS attack or a bug, developers also express a monthly budget that they are willing to pay. XX uses this budget as a guideline and throttles invocations or decreases quality of service in the case that usage is not within reason given the expected budget, though it does not guarantee that the budget will not be exceeded by small amounts.

### 2.2 XX Design

XX has as its goal to enforce the priorities attached to jobs, which means it needs to prefer higher priority jobs over lower ones, and preempt the latter when necessary.

As shown in Figure 1, XX sits behind a load balancer, and consists of: a distributed global scheduler, which places new job invocations, a dispatcher, which runs on each machine and communicates with the global scheduler shards, and a machine scheduler, which enforces priorities on the machines.

Global scheduler shards store the jobs waiting to be placed in a multi queue, with one queue per priority.

**Idle list.** Each global scheduler shard also maintains an *idle list*, which holds machines that have a significant amount of memory available. In the shards idle list each machine's entry is associated with the amount of free memory as well as the current amount of jobs on the machine. The idle list exists because datacenters are large: polling a small number of machines has been shown to be very powerful, but cannot find something that is a very rare occurrence. What this means is that polling is likely to find a machine in the lower quantile of the datacenter, but not at the absolute bottom — it will not find one of the handful of machines that are actually

idle. Having an idle list allows these machines, which are expected to be rare in a high-utilization setting, to make themselves visible to the global scheduler. The idle list also allows the global scheduler to place high priority processes quickly, without incurring the latency overheads of doing the polling to find available resources.

**Machine Scheduler.** The machine scheduler is a preemptive priority scheduler: it preempts lower priority jobs to run higher priority ones. Being unfair and starving low priority jobs is desirable in XX, since image processing jobs should not interrupt a page view request processing, but vice versa is expected. Within priority classes the machine scheduler is first come first served. This matches Linux' 'sched fifo' scheduling.

**Dispatcher.** The dispatcher is in charge of adding itself to a shard's idle list when memory utilization is low. The dispatcher chooses which list to add itself to using power-of-k-choices: it looks at k shards' idle lists and chooses the one with the least other machines in it. If the machine is already on an idle list on shard  $i$ , but the amount of available memory has changed significantly (either by jobs finishing and memory being freed or by memory utilization increasing because of new jobs or memory antagonists), the dispatcher will update shard  $i$ 's idle list. These interactions from the dispatcher to free lists are represented by the blue arrows in Figure 1.

The dispatcher is most importantly in charge of managing memory. When memory pressure occurs, the dispatcher uses priority-based paging to move low priority processes that won't be running soon anyway off the machine's memory. In this scenario, having priority scheduling creates the opportunity that enables this to be realistic: because the dispatcher knows that the lowest priority jobs will not be run until the high priority jobs have all finished, it can page its memory out knowing it will not be needed soon. Doing so, and avoiding a churn of jobs with paged memory, requires that the memory of machines is large against the amount of jobs that would be able to run at once. This assumption ensures that memory pressure high enough to require paging will only occur when there are many more jobs than there are cores, and thus that the paged job will not be running soon. The dispatcher pages the low priority job back in when the memory pressure is gone and the job will be run.

#### Global Scheduler Shards.

Shards choose what job to place next by looking at each job at the head of a queue in the shards multi-queue, and comparing the ratio of priority to amount of time spent in the queue. This ensures that high priority jobs don't have to wait as long as lower priority jobs to be chosen next, but low priority jobs will get placed if they have waited for a while.

When placing the chosen job, the shard finds a machine to run it. The shard will first look in its idle list, and if it is

not empty, will choose the machine with the smallest queue length.

If there are no machines in the idle list, the shard switches over to power-of-k-choices: it polls k machines, getting the amount of jobs running from each. The shard then places the new job on the machine with the smallest queue. It is desirable to have a maximally heterogenous set of priorities on each machine, but since jobs come in randomly the global scheduler does not explicitly have to enforce this.

### 3 Preliminary Results

In order to understand the case for XX, we ask the following questions:

- (1) Are priorities good at meeting SLAs?
- (2) How well do schedulers with no priority information do at meeting SLAs?
- (3) Is XX good at binpacking?

To explore these questions, we built a simulator in go[1], which simulates different scheduling approaches.

#### 3.1 Experimental Setup

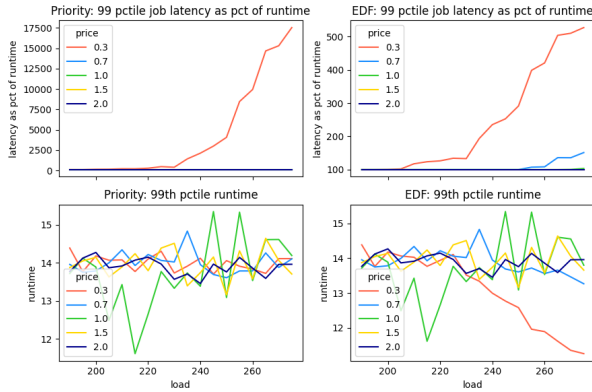
In each version of the simulator, jobs arrive in an open loop at a constant rate. The simulator attaches three main characteristics to each job it generates: runtime, priority, and memory usage. *Job runtime* is chosen by sampling from randomly generated long tailed (in this case pareto) distribution: the relative length of the tail ( $\alpha$  value) remains constant, and the minimum value ( $x_m$ ) is chosen from a normal distribution. This reflects the fact that different functions have different expected runtimes (chosen from a normal distribution), and that actual job runtimes follow long tailed distributions (so each pareto distribution that we sample represents the expected runtime distribution of a given function). *Job priority* is chosen randomly, but weighted: the simulator uses a vaguely bimodal weighting across priorities. The simulator has n different priority values, each assigned to a fictitious price. Because functions are randomly assigned a priority, runtime and priority are not correlated. *Job memory usage* is chosen uniform random between 1MB and 10GB.

When comparing two different simulated schedulers, they each are given an identical workload and then each simulate running that workload.

The simulator makes some simplifying assumptions:

- (1) functions are compute bound, and do not block for i/o
- (2) functions use all of their memory right away
- (3) communication latencies are not simulated

We simulate running 100 machines with 8 cores each, with 4 scheduler shards and a k-choices value of 3 when applicable.



**Figure 2:** a comparison of latency between priority scheduling (left) and EDF (right)

### 3.2 Are priorities good at meeting SLAs?

In the end developers care about job latency, so it is important to understand how well priorities do at reflecting and enforcing SLAs.

In order to define the desired SLA, we assign each job invocation a deadline, and to make things simple we allow the deadline to be a function of the job’s true runtime. Simply setting the deadline to be the runtime is not realistic: perhaps for high priority jobs that is the case, but for low priority jobs that is not true, the reason they are low priority is because they are willing to wait. We thus define the deadline as a function of the runtime as well as the priority, where  $\text{deadline} = \text{runtime} * \text{maxPrice}/\text{price}$  (as if each process were weighted by its price). We simulate an Earliest Deadline First (EDF) scheduler, which is queuing theoretically proven to be optimal in exactly the way we wanted: if it is possible to create a schedule where all jobs meet their deadline, EDF will find it[1].

We compare the latencies observed in this EDF simulated scheduler with an idealized version of a preemptive priority scheduler, and look at the latencies they both get. A good result for XX would show little difference between the two: we expect the EDF version to do better, because it is both theoretically optimal and has access to perfect information, but if the priority scheduler does similarly then we know that it can be a good proxy for deadlines. We can see the results in Figure 2. Although in the EDF version the latency spikes much less, the bottom graph tells us why: EDF, because it has knowledge of jobs’ approximate runtime, stops running the longer low priority jobs.

### 3.3 How well do schedulers without priority info do?

This question is relevant in order to understand if we need priorities at all: is there a scheduler that can, without having any access to information about which jobs are important, still ensure that jobs perform well?

To explore this question, we look at how an existing state of the art research scheduler that does not take any form of priority into account performs on the web server’s workload. We simulate Hermod[1], a state-of-the-art serverless scheduler built specifically for serverless. Hermod’s design is the result of a from-first-principles analysis of different scheduling paradigms. In accordance with the paper’s findings, we simulate least-loaded load balancing over machines found using power-of-k-choices, combined with early binding and Processor Sharing machine-level scheduling. Hermod does not use priorities in its design, and as such the simulator ignores jobs’ priority when simulating Hermod’s design.

We simulate Hermod and XX, and compare the latencies that each achieves, primarily for the page view jobs. A strong result for XX would show that latencies for higher priority jobs start going up at a higher load than for Hermod, which would indicate that Hermod’s scheduling cannot perform as well in high load settings on high priority jobs. Figure ?? shows the results.[hmg: ok but how is this not obvious: we made a new metric and look we did better at it than other systems that did not use that metric wow what a surprise]

### 3.4 Is XX good at binpacking?

Now that we have established that priorities are a good metric for meeting SLAs, the question becomes whether XX’s design allows it to efficiently realize that goal, despite being distributed in nature.

To answer this question, we take an idealized version of the scheduling XX does, and compare XX’s performance to it. The idealized scheduler runs in a centralized setting: it is as if the whole datacenter is one machine with many many cores. As such there is no memory fragmentation, and its utilization represents an optimal solution.[hmg: this feels slightly silly as an experiment, because data points don’t give a context to be able to tell how close is good. Maybe do an ablation study vibe, where we look at our techniques and see which ones help?]

We compare the latencies as well as the memory and compute utilization we observe. A good result for XX would show that the two do not differ too much[hmg: super ill defined]; we expect utilization to have a higher variance for the XX version of the simulator, but a good result would show that we are able to reduce the variance and raise the average and tail utilization by using idle lists. We can see the results in Figure ??.[hmg: TODO do this]

## 4 Discussion

In this section we discuss the the way that having priorities would impact how we can approach some of the open questions that remain about how a scheduler should best manage resources.



#### 4.1 Managing memory

In section 2, XX makes decisions based on maximum memory usage estimates for each job, provided by the developer. However, in reality, making scheduling decisions based off of this sort of estimate of a maximum can lead to under-utilization. Jobs often have highly variable memory usage; for instance the distribution of content among users of social media platforms is notoriously skewed[1], so jobs that fetch information for a user and process it will be run over very different amounts of data depending on who the user is. So making scheduling decisions based off of an estimate of a maximum will lead to memory under-utilization, and instead we are forced to play a game of overprovisioning.

Having priorities can help. We overprovision in order to achieve high utilization, so need a way of dealing with a situation where there are not enough resources for all the jobs on a machine. Under that situation of high memory pressure, we suggest a priority-based paging approach, where the lower priority jobs' memory are paged out. Because machines run highest priority first scheduling, we know that the paged job will not be run until the higher priority jobs have finished (and thus freed their memory).*[hmng: there's a caveat here where they could block on i/o, but the likelihood that all the jobs except the last one block on i/o is low]* Later, when there is lower load, and the previously paged job is ready to run again, the machine can page back in all of its memory. Thus the paged job pays no latency overhead for having been paged, and the amount of time spent paging memory is low: once to move all the memory out, and once to move it back in. This is the minimal amount of data movement necessary to free up the required memory without deleting/killing.

#### 4.2 Managing compute

We have already seen that having priorities helps XX deal with overprovisioning on the compute side. Traditionally, managing the response time for important jobs requires the overall load on the system to be low, because slowdown is averaged across all the jobs that can share resources. However, having priorities allows XX to be targeted about how compute resources are shared: rather than averaging out the slowdown caused by overprovisioning across all the jobs by having them time share the cores, XX can use priorities to decide which job has to wait. This means that in priority scheduling the amount of time a job spends waiting for resources is only defined by the load of jobs with equal or higher priority.

The flipside of this is that it is possible that the entire load of the datacenter will be such that low priority jobs are starved. This is acceptable and in fact desirable for small amounts of time, but keeping this effect in check requires managing the overall load.

We propose to solve this problem by ensuring that it can never be the case that there is so much load on high priority jobs that the data center will be full with them to such a degree that other jobs can't run. There is evidence for and we expect that load on a high level will mostly be stable, with diurnal and annual patterns.[1]

This means that providers can mostly choose the rough breakdown of the load they will have at any given time (ie they can choose a percentage of how many jobs they want in each priority, and change it by adjusting prices or not allowing users to select that level anymore).

### 5 Related Work

Many other projects have explored how to do better serverless scheduling.

Systems like Sparrow[1], Hermod[1], or Kairos[1] improve performance of scheduling in the distributed setting by trying out and using different scheduling policies. Unlike XX, they do not differentiate between individual types of jobs that come in.

Some systems generate information about jobs that are coming in to help placement decisions; for instance ALPS[1], which observes and learns the behaviors of existing functions and then makes scheduling decisions based on those; or Morpheus[1], which learns SLOs from historical runs, and then translates these to recurring reservations. XX instead gets the priorities directly from the developers as part of its interface.

Other existing systems ask developers to specify the desired behavior, like XX does. However, they don't change the actual underlying machine-level scheduling, and expressing priority is more complex than in XX, and/or has a different goal.*[hmng: this feels a bit awkward/not on the nail. Also hard because of unknown unknowns, ie there may be a related work that I don't know about for which this is not true. Hard to make specific enough so that is unlikely, while still actually going beyond just the literal design ]*

Sequoia[1], for instance, creates a metric of QOS for serverless functions. Unlike XX, Sequoia does not implement a new scheduler, but builds a layer in front of AWS lambda.

Another project[1] creates a language of Allocation Priority Policies (APP), which is a declarative language to express policies, then builds a scheduler around that. Unlike XX, the APP language is built around making load balancing decisions and ultimately defines a mapping of jobs to workers, rather than associating priorities with the jobs and having the scheduler do the scheduling automatically.

AWS lambda itself also goes this route, by offering two different ways for developers to influence lambda function scaling: provisioned and reserved concurrency[1]. Provisioned concurrency specifies a number of instances to keep warm, and reserved concurrency does the same but ensures that

those warm instances are kept for a specific function. However, these knobs are more centered around mitigating the effects of cold start rather than actually affecting the way the jobs are scheduled.

## References

- [1] TODO. “TODO”. In: 2020.