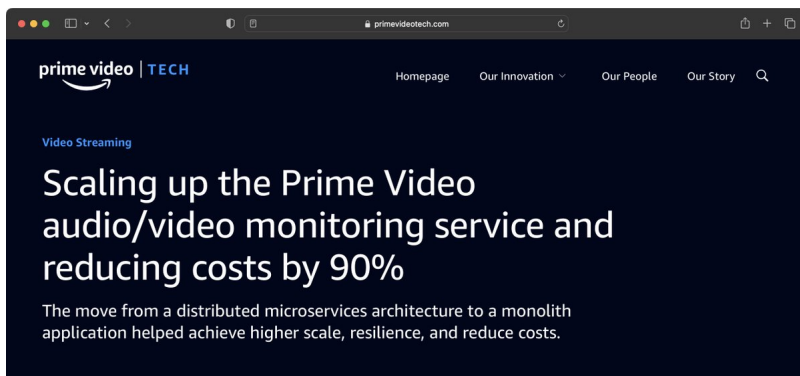


Kelsey Hightower @kelseyhightower *Thu May 04 12:18:56 +0000 2023*

The Amazon Prime Video team was able to reduce cost by moving from Serverless backed by Lambda to monoliths running on VMs.

"Moving our service to a monolith reduced our infrastructure cost by over 90%. It also increased our scaling capabilities."

<https://t.co/lnIWWveqVj> <https://t.co/qt4lnonN3p>



Marcin Kolny
Mar 22, 2023

At Prime Video, we offer thousands of live streams to our customers. To ensure that customers seamlessly receive content, Prime Video set up a tool to monitor every stream viewed by customers. This tool allows us to automatically identify perceptual quality issues (for example, block corruption or audio/video sync problems) and trigger a process to fix them.

Our Video Quality Analysis (VQA) team at Prime Video already owned a tool for audio/video quality inspection, but we never intended nor designed it to run at high scale (our target was

Most popular

"We're just beginning to build the future of live sports streaming"

Feb 07, 2023

Prime Video announces Amazon Research Awards recipients for fall 2022

This isn't a dig against Lambda as that platform helped the team build the service fast and get to market.

"We designed our initial solution as a distributed system using serverless components, which was a good choice for building the service quickly."

But it is a testament to the overhead of microservices in the real world. Moving data around is typically an underestimated cost.

"The second cost problem we discovered was about the way we were passing video frames (images) around different components."

A monolithic architecture doesn't mean a spaghetti code base. You should be writing modular code regardless of the deployment model.

"Conceptually, the high-level architecture remained the same. We still have exactly the same components as we had in the initial design."

The post isn't about microservices vs monoliths. It's about using the right tool for the job.

"Microservices and serverless components are tools that do work at high scale, but whether to use them over monolith has to be made on a case-by-case basis."