Chair for Decentralized Information Systems and Data Management
TUM School of Computation, Information and Technology
Technical University of Munich

# Cloud Information Systems

## Exercise 9

16th December 2024

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter
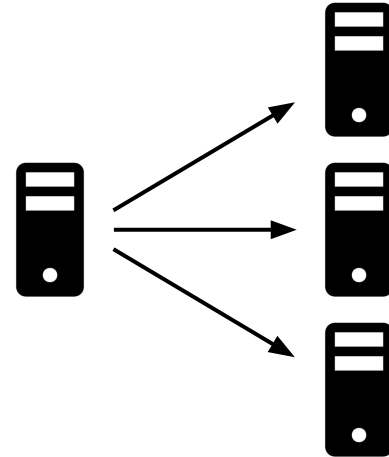
# 1. Recap: Scalability

vertical scaling ("scale up")
→ bigger machines

horizontal scaling ("scale out")
→ more machines

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter | Chair for Decentralized Information Systems and Data Management | Technical University of Munich

2

# 1. Recap: Scalability

| Cloud makes scale up easier | Cloud favors scale out architecture |
|---|---|
| - migrate to bigger VM slice<br>- live migration would make this quite easy<br>- enables high-bandwidth and low-latency communication<br><br><br>but: very large servers can become uneconomical | - add/remove machines when workload changes<br>- enables fault tolerance through redundancy<br>- works nicely for stateless services (eg, web servers)<br><br>but: downsides like<br>    → distributed state management,<br>    → new failure modes,<br>    → bottleneck: network bandwidth and latency |

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter | Chair for Decentralized Information Systems and Data Management | Technical University of Munich
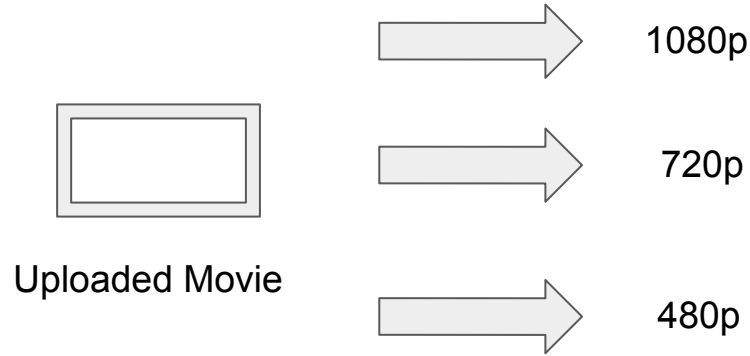
3

# 1. Recap: Scalability

If the problem is not trivially partitionable, it is best to first

- scale up to one medium-size server and
- only if that is not enough to scale out to a cluster of medium-size servers

+ exploits fast communication
+ hardware resources are generally proportional to cost
+ enables "arbitrary" scaling
- this implies one has to exploit two levels of parallelism

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter | Chair for Decentralized Information Systems and Data Management | Technical University of Munich

4

# 2. Live Exercise: Netflix Video Transcoding

"While onboarding more streams to the service, we noticed that running the infrastructure at a high scale was very expensive."

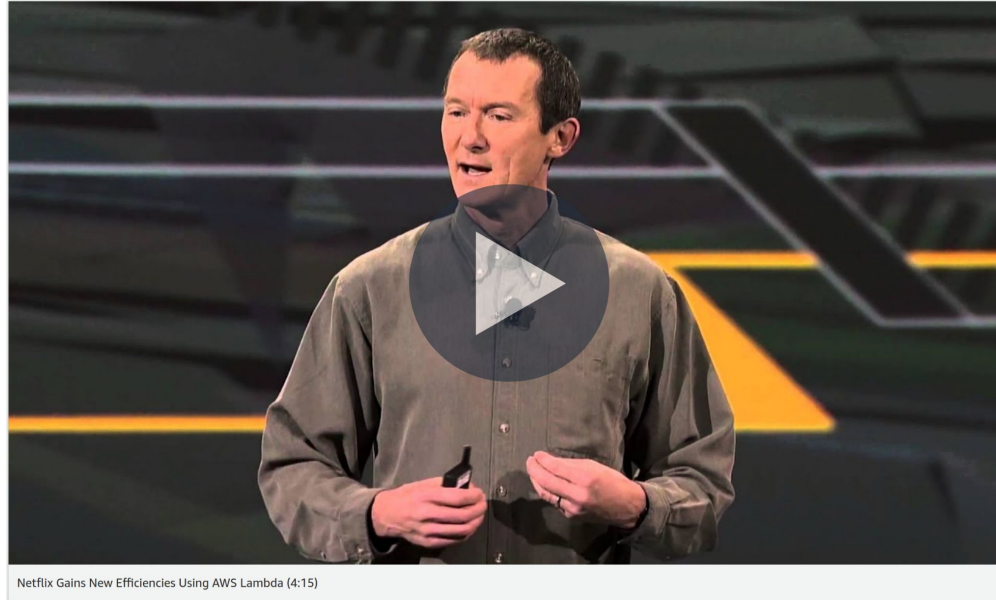Scaling up the Prime Video audio/video monitoring service and reducing costs by 90%

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter | Chair for Decentralized Information Systems and Data Management | Technical University of Munich

5

# 2. Live Exercise: Netflix Video Transcoding

1080p

720p

Uploaded Movie

480p

➔ Publisher Uploads new Movie to Netflix
➔ Different Customers have different Demands (4K TV vs Mobile Phone) and Subscription Models (Basic vs Standard vs Premium)

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter | Chair for Decentralized Information Systems and Data Management | Technical University of Munich

6

# 2. Live Exercise: Netflix Video Transcoding

Netflix Gains New Efficiencies Using AWS Lambda (4:15)

Netflix Keynote

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter | Chair for Decentralized Information Systems and Data Management | Technical University of Munich

7

# 2. Live Exercise: Netflix Serverless Video Transcoding



Split Movie into Segments of 5 Minutes

Perform Transcoding (1 Lambda / Segment / Encoding)

*new video event*

*new segment event*

*new segment event*

*new segment event*

upload

**S3 bucket**

**serverless segmenter**

**S3 bucket**

**serverless transcoder**

**serverless transcoder**

**serverless transcoder**

**S3 bucket**

For Details, see: Comer, "The Cloud Computing Book"

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter | Chair for Decentralized Information Systems and Data Management | Technical University of Munich

8

# Q: What is the Cost of Transcoding a two hour Movie ?

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter | Chair for Decentralized Information Systems and Data Management | Technical University of Munich

9

# 2. Live Exercise: Assumptions

- **General Parameters**:
  - 2 hour Movie in 4K Quality
  - 30 Frames per Second
  - 5 Minutes per Segment
  - 6 Output Resolutions (Assumption based on # of Video Qualities on YouTube)
- **Video Transcoding:**
  - We use FFMPEG on Graviton Instances
  - **c6g.4xlarge:** 41 FPS
  - **c7g.4xlarge:** 59 FPS

  AWS Graviton Benchmark

- Frames are independent and can be processed in parallel -> Processed FPS scale linearly with the number of vCPU cores utilized
- We can ignore S3 Storage Cost -> Actual Deployment happens via CDN

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter | Chair for Decentralized Information Systems and Data Management | Technical University of Munich

10

# 3. Questions

Prof. Dr. Viktor Leis, M.Sc. Till Steinert, M.Sc. Jana Vatter | Chair for Decentralized Information Systems and Data Management | Technical University of Munich

11