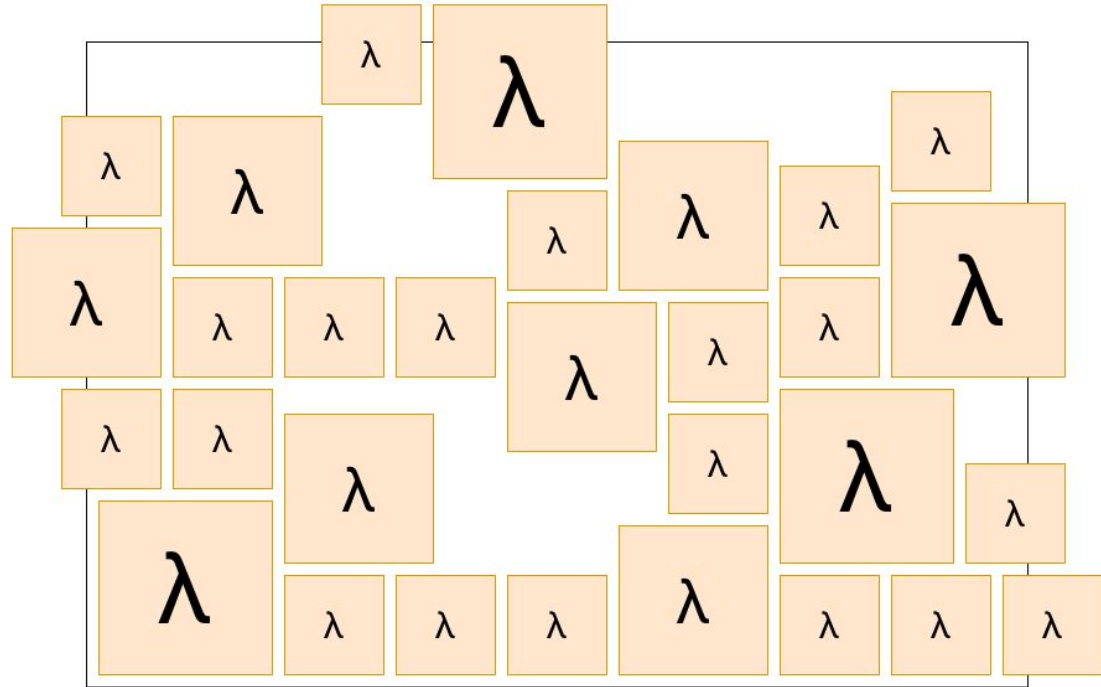
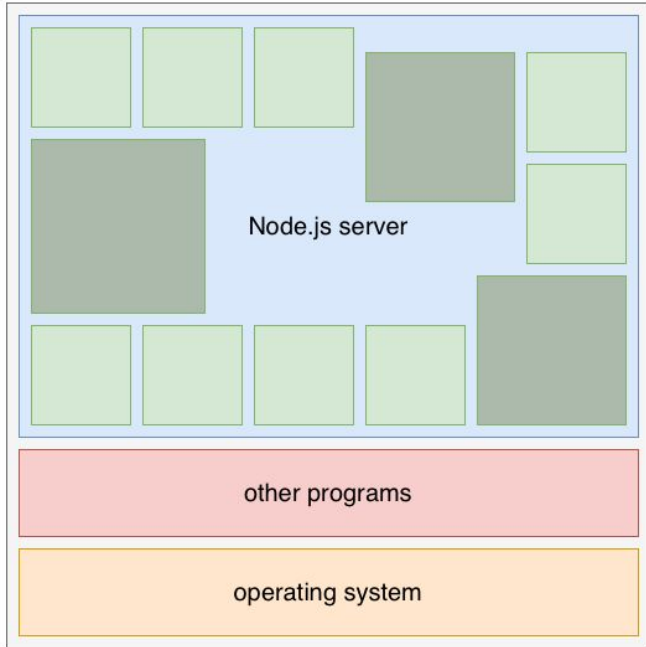


# AWS Serverless: Scaling small and large applications

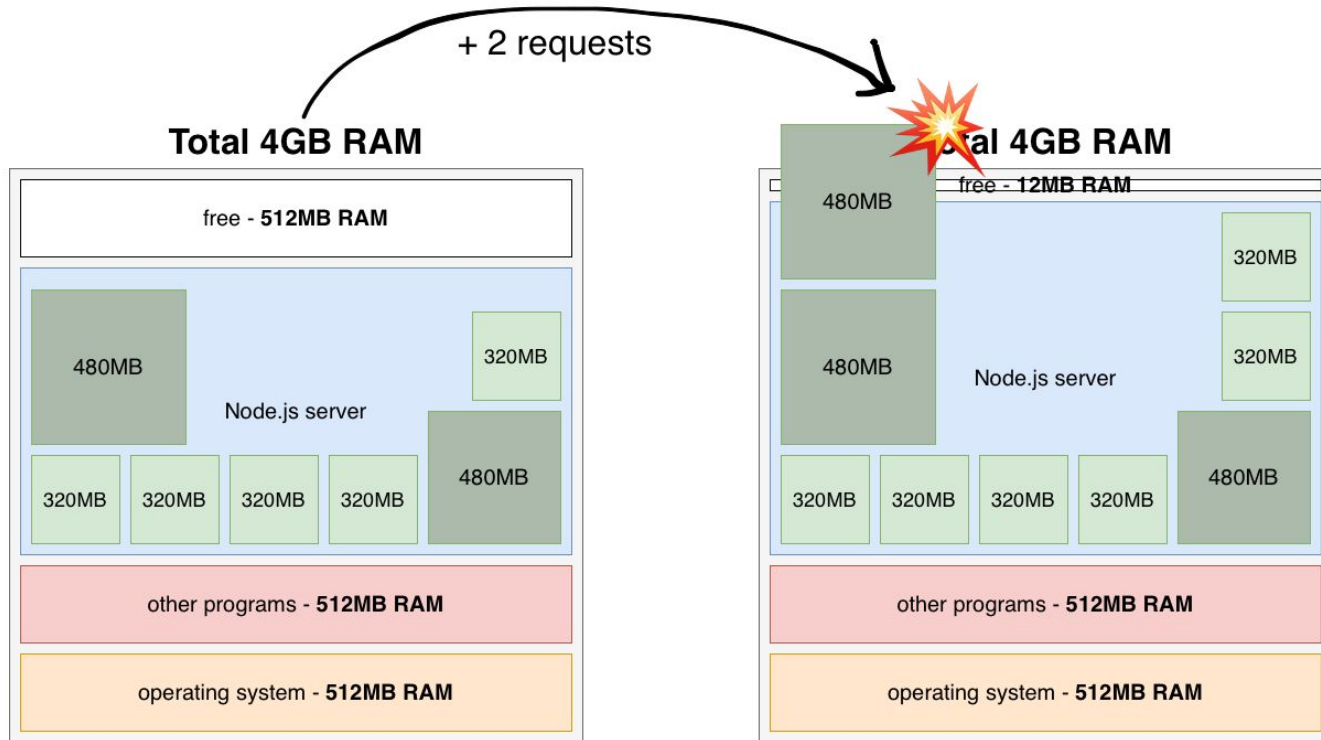
👋 Show of hands



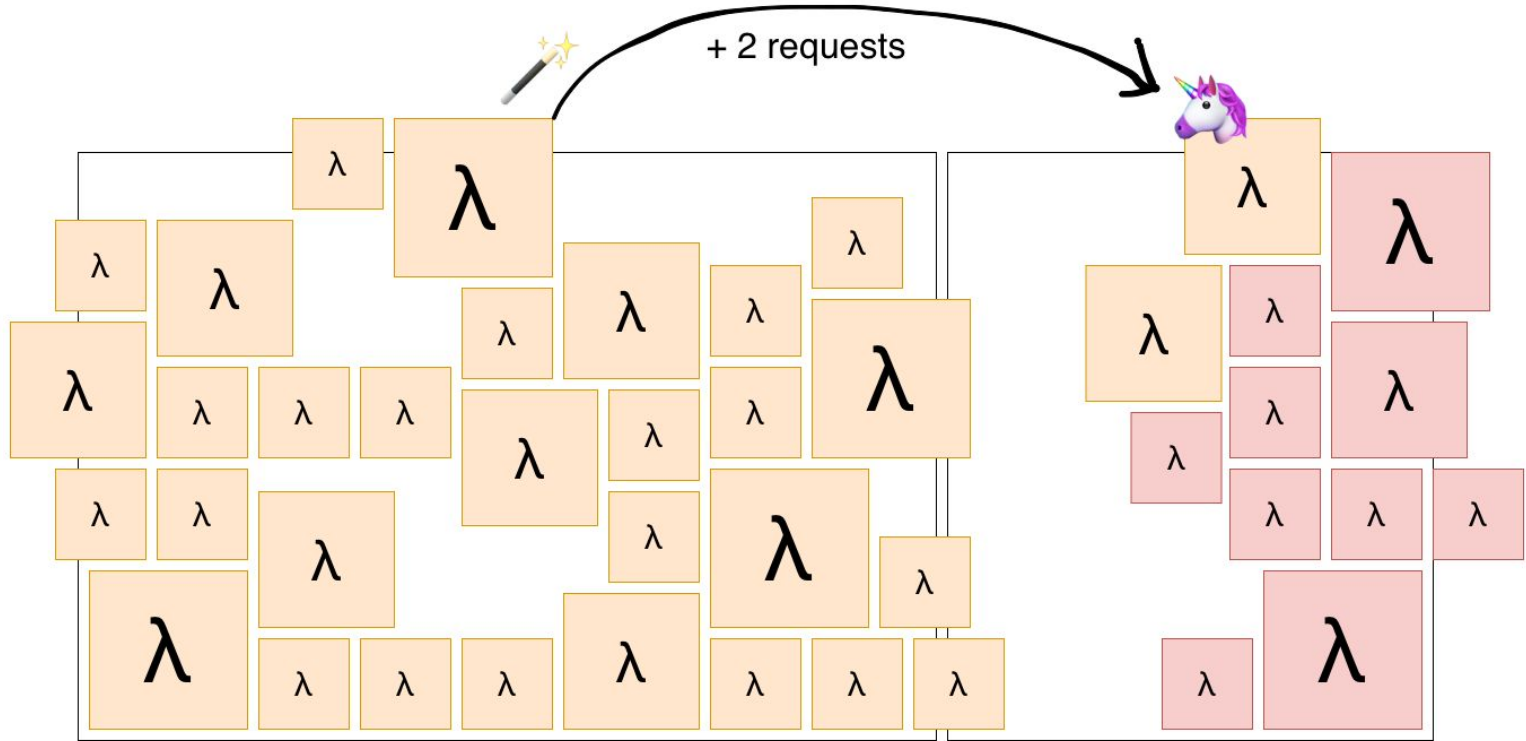
# Traditional VM vs Lambda



# Traditional VM scaling



# Lambda scaling

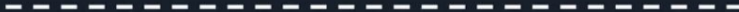




# The biggest secret of Serverless



(and many other serverless services)



Hundreds of thousands of EC2 instances under-the-hood



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Usman Khalid

He/His/Him

Head of Serverless Compute

[linkedin.com/in/ahmedusmankhalid/](https://www.linkedin.com/in/ahmedusmankhalid/)

hundreds of thousands of servers.



# Lambda limits

15m hard timeout

128 MB - 10240 MB RAM

1 - 6 vCPU

cold starts

telemetry (ADOT has its own limitations)

limited official runtimes

isolation - public service

(not great for highly regulated subjects)



# Breaking the limits, slowly and steady

15m hard timeout	Durable functions
128 MB - 10240 MB RAM	Lambda Managed Instances (LMI)
1 - 6 vCPU	up to 32 GB RAM + 16 vCPU
cold starts	SnapStart (Java, Python, .Net)
telemetry (ADOT has its own limitations)	OpenTelemetry (Otel) Support <a href="#">#20</a>
limited official runtimes	Container Images, <a href="#">Lambda Web Adapter</a>
isolation - public service (not great for highly regulated subjects)	LMI - runs inside your VPC (not same as VPC support - accessing services in VPC)

# AWS Lambda — A Decade+ of Serverless Innovation

43 milestones · Nov 2014 → Apr 2026

2014	2016	2018	2019	2021	2023	2025	2026
<p>NOV 13, 2014</p> <p><b>AWS Lambda Announced</b></p> <p>Unveiled at re:Invent — Node.js runtime, S3 / DynamoDB / Kinesis event triggers.</p>	<p>FEB 11, 2016</p> <p><b>VPC Access</b></p> <p>Functions can reach private VPC resources.</p>	<p>JAN 15, 2018</p> <p><b>Go + .NET 2.0</b></p> <p>Go and .NET Core 2.0 added as managed runtimes.</p>	<p>SEP 2019</p> <p><b>Hyperplane VPC Networking</b></p> <p>Hyperplane ENIs deliver fast cold starts inside VPCs.</p>	<p>MAY 24, 2021</p> <p><b>Extensions GA</b></p> <p>Lambda Extensions API GA; observability/security partners.</p>	<p>APR 7, 2023</p> <p><b>Response Streaming</b></p> <p>Stream response payloads incrementally — over 6 MB.</p>	<p>2025</p> <p><b>SnapStart on arm64</b></p> <p>SnapStart on Graviton2 arm64; additional regions added.</p>	<p>APR 2026</p> <p><b>S3 Files</b></p> <p>Mount S3 buckets as file systems inside Lambda — file API on S3 durability.</p>
<p><b>2015</b></p> <p>APR 9, 2015</p> <p><b>General Availability</b></p> <p>Production-ready event-driven compute service.</p>	<p>NOV 18, 2016</p> <p><b>SAM + Env Variables</b></p> <p>Serverless Application Model; environment variables.</p>	<p>JUN 28, 2018</p> <p><b>SQS Event Source</b></p> <p>Amazon SQS as a Lambda event source.</p>	<p>NOV 25, 2019</p> <p><b>Async Destinations</b></p> <p>Destinations + retry and event-age controls for async.</p>	<p>SEP 29, 2021</p> <p><b>Graviton2 (arm64)</b></p> <p>arm64 (Graviton2) support — better price / performance.</p>	<p>AUG 2023</p> <p><b>12× Faster Scaling</b></p> <p>1,000 concurrent executions every 10 sec per function.</p>	<p>AUG 2025</p> <p><b>INIT Phase Billing</b></p> <p>Standardized billing for the function initialization phase.</p>	
<p>APR 9, 2015</p> <p><b>SNS Event Source</b></p> <p>SNS notifications can trigger Lambda functions.</p>	<p>DEC 1, 2016</p> <p><b>C# + Lambda@Edge</b></p> <p>.NET Core runtime; Lambda@Edge runs at CloudFront edge.</p>	<p>OCT 10, 2018</p> <p><b>15-minute Timeout</b></p> <p>Maximum execution duration raised from 5 to 15 minutes.</p>	<p>DEC 3, 2019</p> <p><b>Provisioned Concurrency</b></p> <p>Pre-warmed execution environments — no cold starts.</p>	<p><b>2022</b></p> <p>MAR 23, 2022</p> <p><b>10 GB Ephemeral /tmp</b></p> <p>/tmp storage configurable up to 10 GB per invocation.</p>	<p>NOV 9, 2023</p> <p><b>Amazon Linux 2023</b></p> <p>AL2023 base for managed runtimes and containers.</p>	<p>DEC 2025</p> <p><b>Durable Functions</b></p> <p>Multi-step apps and AI flows; suspend up to 1 year, checkpoint.</p>	
<p>JUN 15, 2015</p> <p><b>Java Runtime</b></p> <p>Java 8 added as a managed runtime.</p>	<p>NOV 29, 2018</p> <p><b>Layers + Custom Runtimes</b></p> <p>Reusable layers plus custom runtimes — bring any language.</p>	<p>NOV 29, 2018</p> <p><b>Ruby + ALB Integration</b></p> <p>Ruby runtime; Application Load Balancer target.</p>	<p>DEC 1, 2020</p> <p><b>Container Images</b></p> <p>Package functions as OCI container images, up to 10 GB.</p>		<p>APR 6, 2022</p> <p><b>Function URLs</b></p> <p>Dedicated HTTPS endpoints for Lambda functions.</p>	<p>2024</p> <p><b>Recursive Loop Detection</b></p> <p>Auto-detect and stop runaway recursive invocations.</p>	<p>DEC 2025</p> <p><b>Managed Instances</b></p> <p>Capacity providers for GPU and high-memory workloads.</p>
<p>APR 9, 2015</p> <p><b>API Gateway</b></p> <p>HTTP invocation via Amazon API Gateway.</p>	<p>NOV 28, 2017</p> <p><b>Traffic Shifting</b></p> <p>Aliases enable canary and linear deployments.</p>	<p>NOV 29, 2018</p> <p><b>Ruby + ALB Integration</b></p> <p>Ruby runtime; Application Load Balancer target.</p>	<p>DEC 1, 2020</p> <p><b>1ms Billing Granularity</b></p> <p>Billing rounded down to the nearest millisecond.</p>	<p>NOV 28, 2022</p> <p><b>SnapStart for Java</b></p> <p>Sub-second cold starts for Java functions.</p>	<p>2024</p> <p><b>VS Code Console Editor</b></p> <p>VS Code-based editor inside the AWS Lambda console.</p>	<p>DEC 2025</p> <p><b>Tenant Isolation</b></p> <p>Separate execution environment per tenant or end-user.</p>	
<p>OCT 8, 2015</p> <p><b>Python + Versioning</b></p> <p>Python 2.7 runtime, versioning, 5-minute execution timeout.</p>	<p>NOV 28, 2017</p> <p><b>Traffic Shifting</b></p> <p>Aliases enable canary and linear deployments.</p>	<p>NOV 29, 2018</p> <p><b>Ruby + ALB Integration</b></p> <p>Ruby runtime; Application Load Balancer target.</p>	<p>DEC 1, 2020</p> <p><b>10 GB Memory / 6 vCPUs</b></p> <p>Memory raised to 10 GB, up to 6 vCPUs per function.</p>	<p>NOV 28, 2022</p> <p><b>SnapStart for Java</b></p> <p>Sub-second cold starts for Java functions.</p>	<p>NOV 2024</p> <p><b>SnapStart for Python / .NET</b></p> <p>SnapStart extended to Python 3.12+ and .NET 8+ runtimes.</p>	<p>DEC 2025</p> <p><b>Kafka + SQS Updates</b></p> <p>Native Kafka schema evolution; SQS provisioned mode.</p>	
	<p><b>2017</b></p> <p>APR 19, 2017</p> <p><b>AWS X-Ray Tracing</b></p> <p>Distributed tracing for Lambda functions.</p>	<p>NOV 29, 2018</p> <p><b>Ruby + ALB Integration</b></p> <p>Ruby runtime; Application Load Balancer target.</p>	<p>DEC 1, 2020</p> <p><b>1 MB Billing Granularity</b></p> <p>Billing rounded down to the nearest millisecond.</p>	<p>NOV 28, 2022</p> <p><b>SnapStart for Java</b></p> <p>Sub-second cold starts for Java functions.</p>	<p>NOV 2024</p> <p><b>SnapStart for Python / .NET</b></p> <p>SnapStart extended to Python 3.12+ and .NET 8+ runtimes.</p>	<p>DEC 2025</p> <p><b>1 MB Async Payload</b></p> <p>Max payload 256 KB → 1 MB for async, SQS, EventBridge.</p>	
	<p>APR 19, 2017</p> <p><b>AWS X-Ray Tracing</b></p> <p>Distributed tracing for Lambda functions.</p>	<p>NOV 29, 2018</p> <p><b>Ruby + ALB Integration</b></p> <p>Ruby runtime; Application Load Balancer target.</p>	<p>DEC 1, 2020</p> <p><b>10 GB Memory / 6 vCPUs</b></p> <p>Memory raised to 10 GB, up to 6 vCPUs per function.</p>	<p>NOV 28, 2022</p> <p><b>SnapStart for Java</b></p> <p>Sub-second cold starts for Java functions.</p>	<p>NOV 2024</p> <p><b>SnapStart for Python / .NET</b></p> <p>SnapStart extended to Python 3.12+ and .NET 8+ runtimes.</p>	<p>DEC 2025</p> <p><b>Console-to-IDE</b></p> <p>One-click IDE handoff; LocalStack offline testing.</p>	

■ Foundational

■ Runtimes

■ Integrations / Networking

■ Performance / Cold-start

■ Architecture / Compute

■ Developer Experience

■ Limits / Pricing




# Kiro + Powers

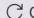


## AWS Observability

Comprehensive AWS observability platform combining CloudWatch Logs, Metrics, Alarms, Ap...  
by AWS

 Try power


 Uninstall

 Check for updates

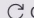


## IAM Policy Autopilot

AWS IAM Policy Autopilot analyzes your application code locally to generate and deploy identity-based p...  
by AWS

 Try power


 Uninstall

 Check for updates




## Build applications with Aurora DSQL

Build applications using a serverless, PostgreSQL-compatible database with scale-to-zero and pay-per-...  
by AWS

 Try power


 Uninstall

 Check for updates

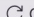


## Build AWS infrastructure with CDK and CloudFormation

Build well-architected AWS infrastructure with CDK using latest documentation, best practices and cod...  
by AWS

 Try power


 Uninstall

 Check for updates

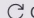


## AWS SAM

An opinionated Kiro Power to aid development with AWS Serverless Application Model (SAM). Includes ...  
by AWS

 Try power


 Uninstall

 Check for updates

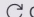


## Build applications with Lambda durable functions

Build resilient multi-step applications and AI workflows with automatic state persistence, retry logic, and ...  
by AWS

 Try power

 Uninstall

 Check for updates

<https://kiro.dev>

# Lambda is not the only serverless compute option

## ECS Fargate

- select # of vCPU and memory
- upload your Docker image
- quick scaling
- fully managed

## ECS Managed Instances

- select the instance type
- managed - no maintenance
- 15% maintenance fee

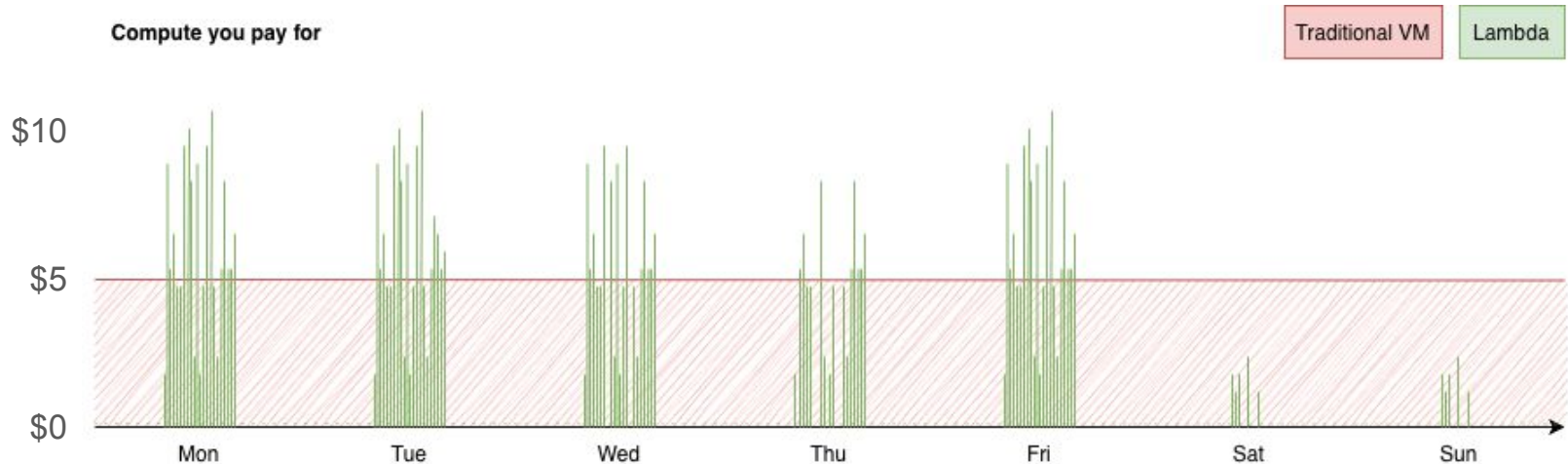


# Traditional VM price comparison

Provider	vCPU / RAM / SSD	Price (original in <b>bold</b> )
Hetzner	2 / 4GB / 40GB (ARM chip)	\$5.85/month <b>€4.99/month</b>
Fly.io	2 / 4GB / 40GB (SSD?)	<b>\$32.68/month</b> €27.88/month
DigitalOcean	2 / 4GB / 80GB	<b>\$24.00/month</b> €20.47/month
Scaleway	2 / 4GB / N/A	\$19.68/month <b>€16.79/month</b>
AWS t4g.medium (70% utilization)	2 / 4GB / 40GB (ARM chip)	<b>\$22.29/month</b> €19.01/month

# Traditional VM limits

- scalability - painful to set up
- durability - single box
- maintenance cost
- overprovisioning - peak load provisioning



# Premature optimization trap

*1 GB of RAM per request?*

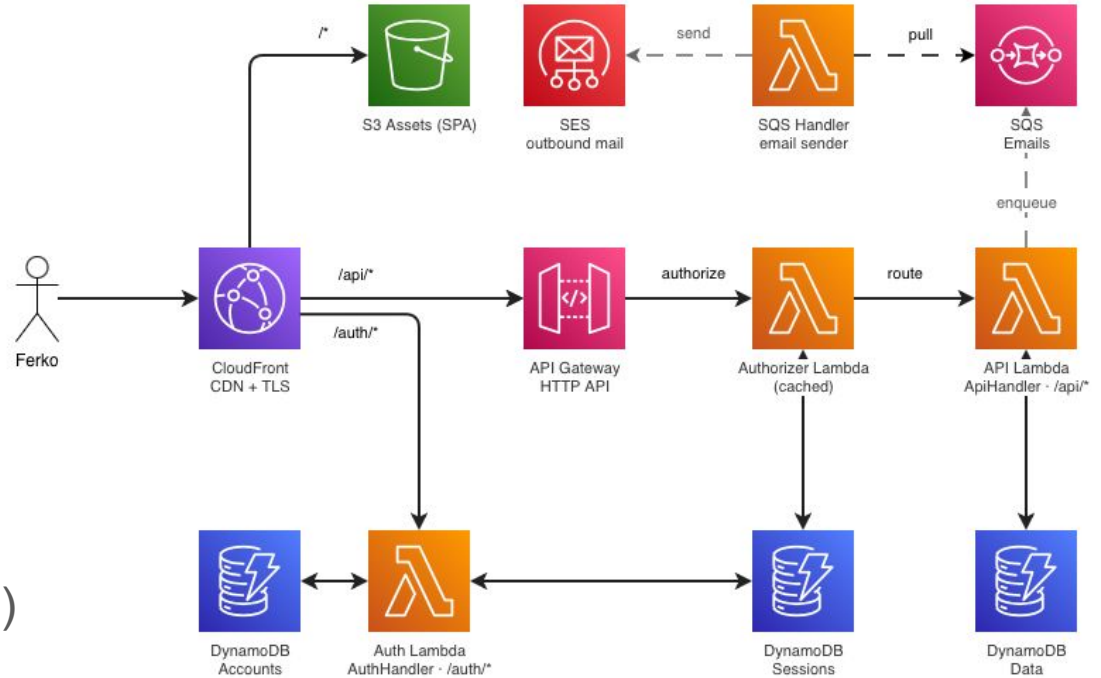
*You can handle X requests with Rust/Golang on single Raspberry Pi!*

## **Instead:**

- set up billing alerts
- set up budgets for Lambda/DynamoDB
- consider cost of infrastructure per paying customer

# Reporting system - simplified architecture

- quarterly reports
- 600 active users
- fully serverless
  
- TanStack router + MUI
- tRPC
- valibot + formisch
- AWS SDK + CDK
- react-email (pre-rendered)
- Node.js 24



# Cost over time (last 5 months)

## Cost and usage overview [info](#)

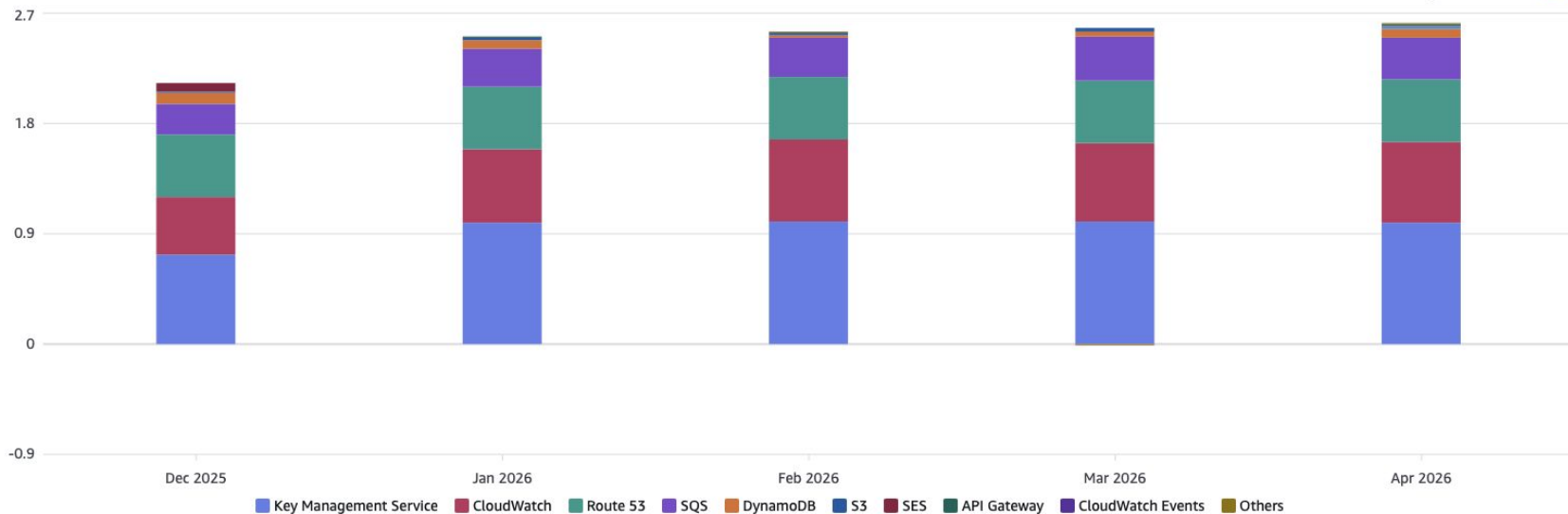
Total cost  
\$12.39

Average monthly cost  
\$2.48

Service count  
19

## Cost and usage graph

Costs (\$)



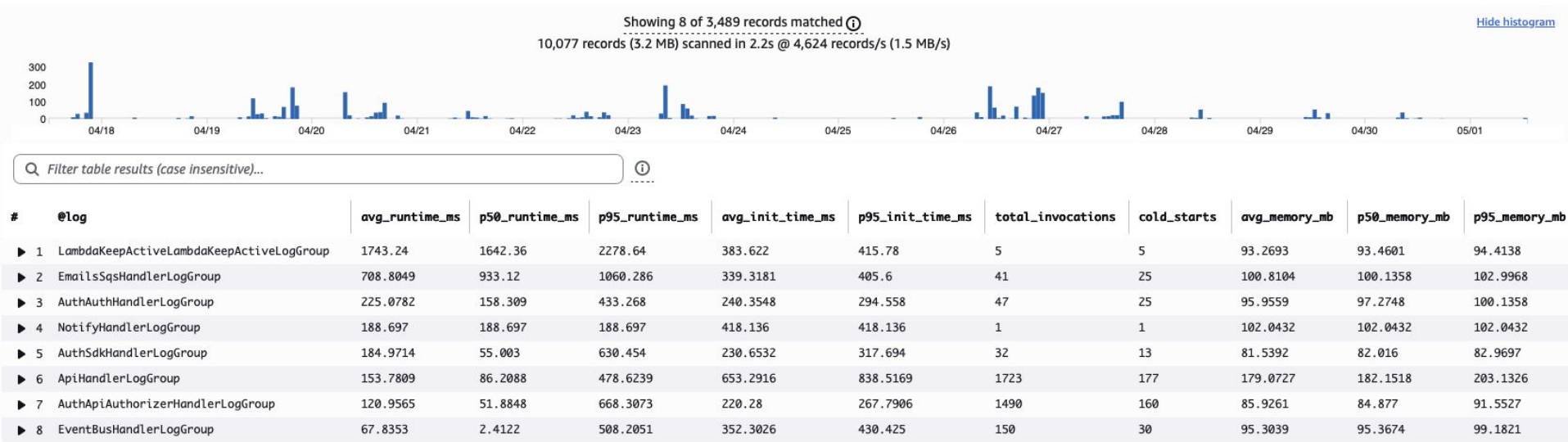
# April 2026 bill breakdown

- **KMS** - used for DNSSec
- **Route 53** - \$0.50 per hosted zone
- **SQS** - even empty queues are being polled constantly
- **S3** - ~500 write/5.5k read reqs
- **API Gateway** - 12k requests
- **CloudFront** - 22k reqs
- **SES** - 300 emails sent

Total active services	Total pre-tax service charges in USD
<b>16</b>	<b>USD 2.62</b>
<input checked="" type="checkbox"/> Key Management Service	USD 0.99
<input checked="" type="checkbox"/> CloudWatch	USD 0.66
<input checked="" type="checkbox"/> Route 53	USD 0.51
<input checked="" type="checkbox"/> Simple Queue Service	USD 0.34
<input checked="" type="checkbox"/> DynamoDB	USD 0.07
<input checked="" type="checkbox"/> Simple Storage Service	USD 0.03
<input checked="" type="checkbox"/> API Gateway	USD 0.01
<input checked="" type="checkbox"/> CloudWatch Events	USD 0.00
<input checked="" type="checkbox"/> CloudFront	USD 0.00
<input checked="" type="checkbox"/> Data Transfer	USD 0.00
<input checked="" type="checkbox"/> CloudFormation	USD 0.00
<input checked="" type="checkbox"/> CloudTrail	USD 0.00
<input checked="" type="checkbox"/> Glue	USD 0.00
<input type="checkbox"/> Lambda	USD 0.00
└─ <input type="checkbox"/> EU (Frankfurt)	USD 0.00
└─ <input type="checkbox"/> AWS Lambda EUC1-Lambda-GB-Secor	USD 0.00
└─ AWS Lambda - Compute Free Tier - 40 / 1,963.037 Lambda-GB-Second	USD 0.00
└─ <input type="checkbox"/> AWS Lambda EUC1-Lambda-GB-Secor	USD 0.00
└─ AWS Lambda - Compute Free Tier for / 2,320.219 Lambda-GB-Second	USD 0.00
└─ <input type="checkbox"/> AWS Lambda EUC1-Request	USD 0.00
└─ AWS Lambda - Requests Free Tier - 1,C / 7 Request	USD 0.00
└─ <input type="checkbox"/> AWS Lambda EUC1-Request-ARM	USD 0.00
└─ AWS Lambda - Requests Free Tier for / 22,803 Requests	USD 0.00
<input checked="" type="checkbox"/> Simple Email Service	USD 0.00
<input checked="" type="checkbox"/> Simple Notification Service	USD 0.00

# Lambda logs breakdown (last 2 weeks)

- runtime in ms
- cold starts
- total invocations
- memory usage in MB



# DynamoDB

- manage your schema in application layer (with `valibot`)
- add schema version to each item for easier management
- filter/sort data in memory (for small datasets)
- implement counters/aggregations with transactions/DynamoDB streams
- combine with OpenSearch/Athena when needed

[Cost-Aware Modeling in DynamoDB:  
Designing for Performance & Spend](#)

by Alex DeBrie



# DSQL

- distributed PostgreSQL-compatible database
- same principles for distributed systems as DynamoDB applies
- highly available and scalable
- fully serverless
- compute on demand + durable storage

## **not a PostgreSQL drop-in replacement**

- missing foreign keys
- everything must be transaction (optimistic concurrency control)
- async index creation

# Serverless tips



presented at  
AWS Community Day Slovakia 2026



# Be honest with yourself

What is the actual scale of the project?

Are you going to deploy internationally or regional?

How many users are you going to serve?

Do you really need relational database?

Are you building for the purpose or out of habit?

Are you the only contributor?



# Ivan Barlog

AWS Solutions Architect



Github [ivanbarlog](#)  
[beesolve](#)

Email [ivan@barlog.sk](mailto:ivan@barlog.sk)

Web [barlog.sk](#)  
[beesolve.com](#)

LinkedIn [ivan-barlog](#)

# Go build something!