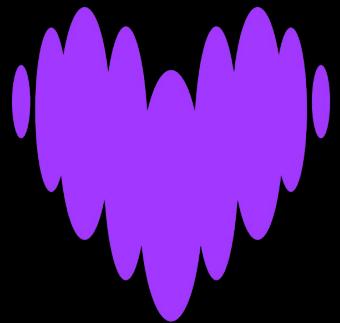


# The journey from (not) REST to GraphQL at scale

GraphQL meetup 22-10-2025





# DEEZER

~10M subscribers

**550 employees** spread over France,  
Germany, United Kingdom, Brazil and USA

Deezer is available in **180+ countries** and **26 languages**

~**800 servers** in Paris + Google Cloud

~**200 network equipments** (switches, routers, etc.)

More than **6 PB** of storage needs

**6 To** of daily logs and **45M** active metrics series

~**4B** requests per day



Desktop



Mobile



TV



Car



Connected  
Speakers



Wearables

~>whoami\_



**Loïc Doubinine**

Backend Developer

@Deezer since 2016



@ztec.fr



@ztec@mamot.fr



@ztec6



[www.ztec.fr/en/social](http://www.ztec.fr/en/social)



[blog.ztec.fr](http://blog.ztec.fr)

# The journey from (not) REST to GraphQL at scale

GraphQL meetup 22-10-2025



(r



THIS IS A TRUE STORY.

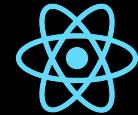
DEEZER



**Context**

# Modernizing old APIs

IOS      ANDROID      WEB      OTHERS



WEB

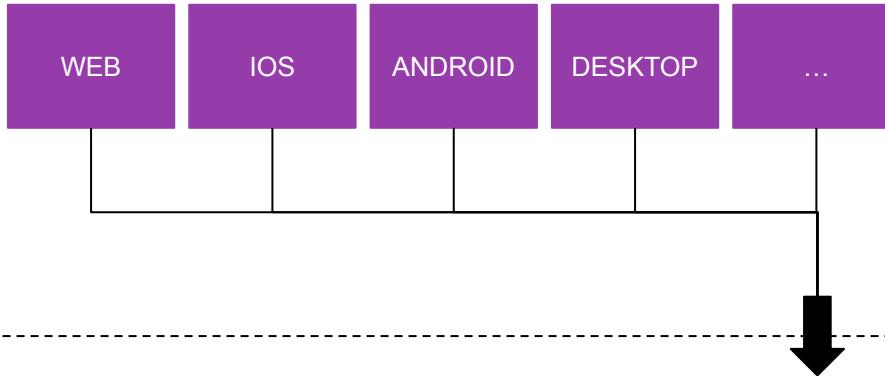
IOS

ANDROID

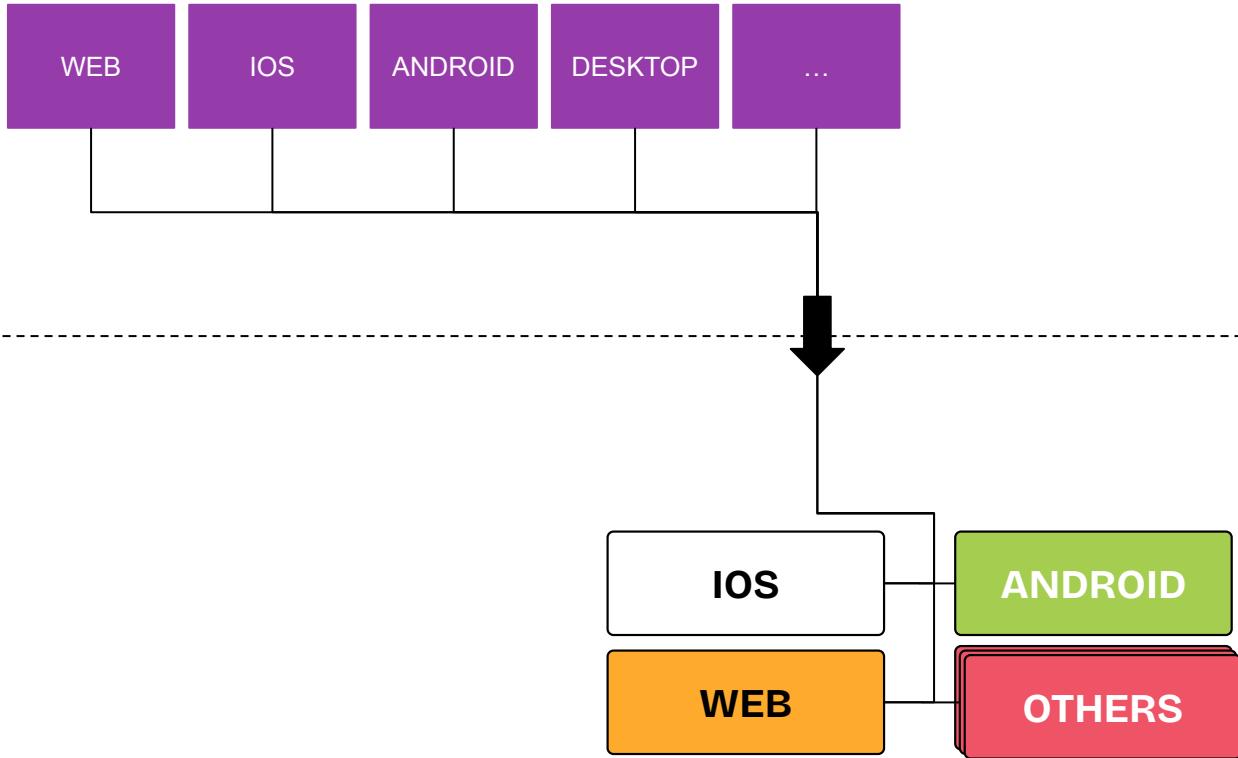
DESKTOP

...

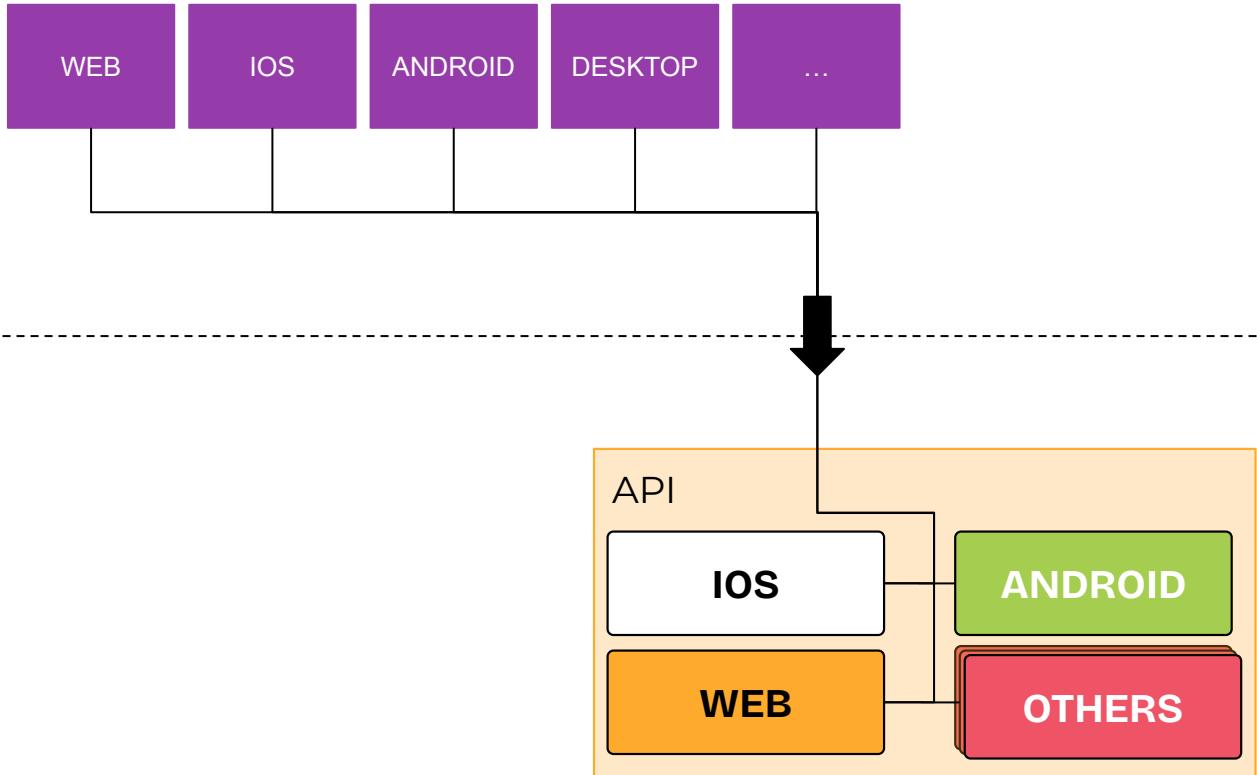




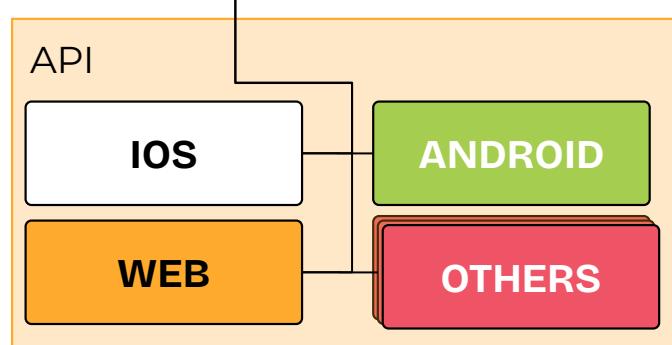
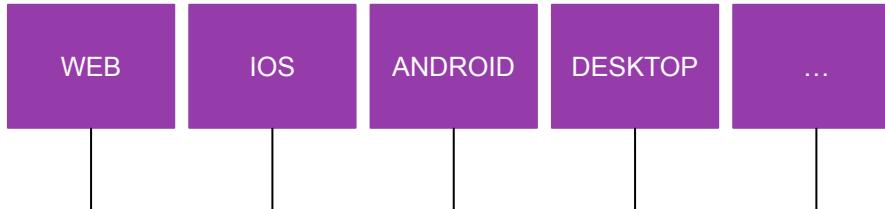
Client & Devices



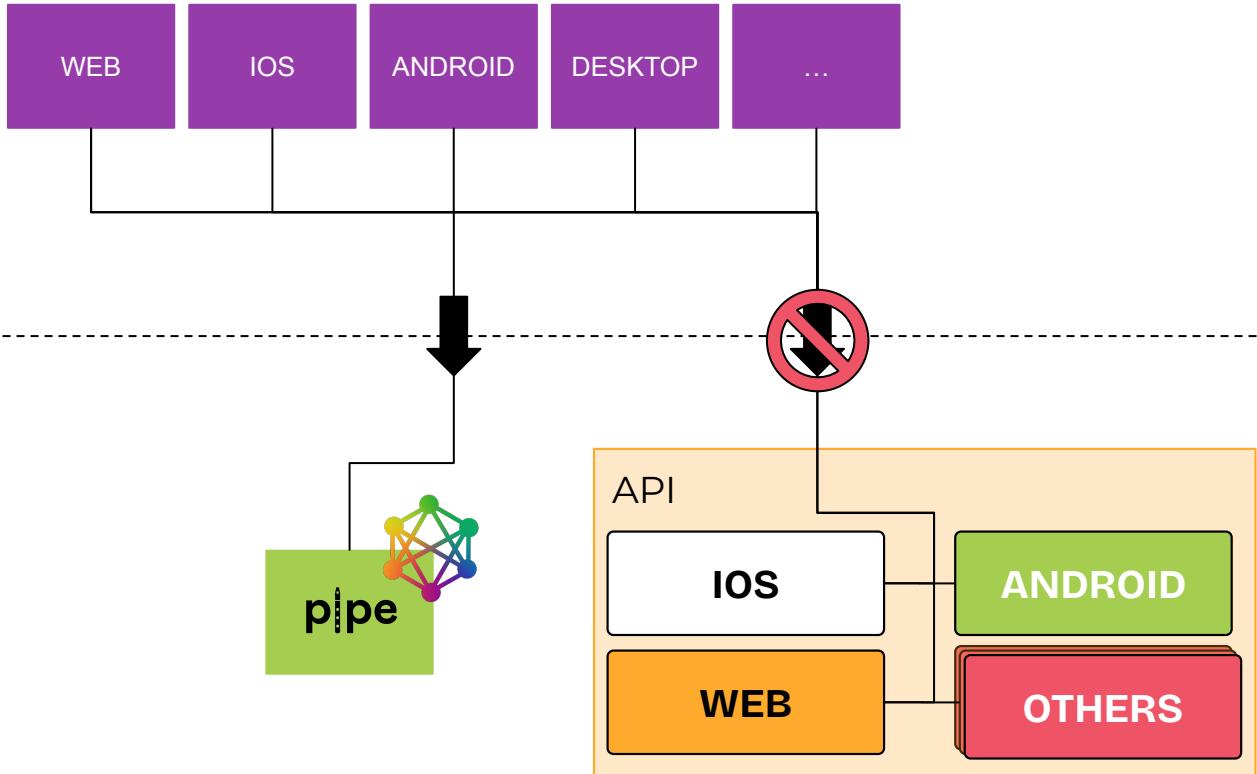
internal  
Client & Devices



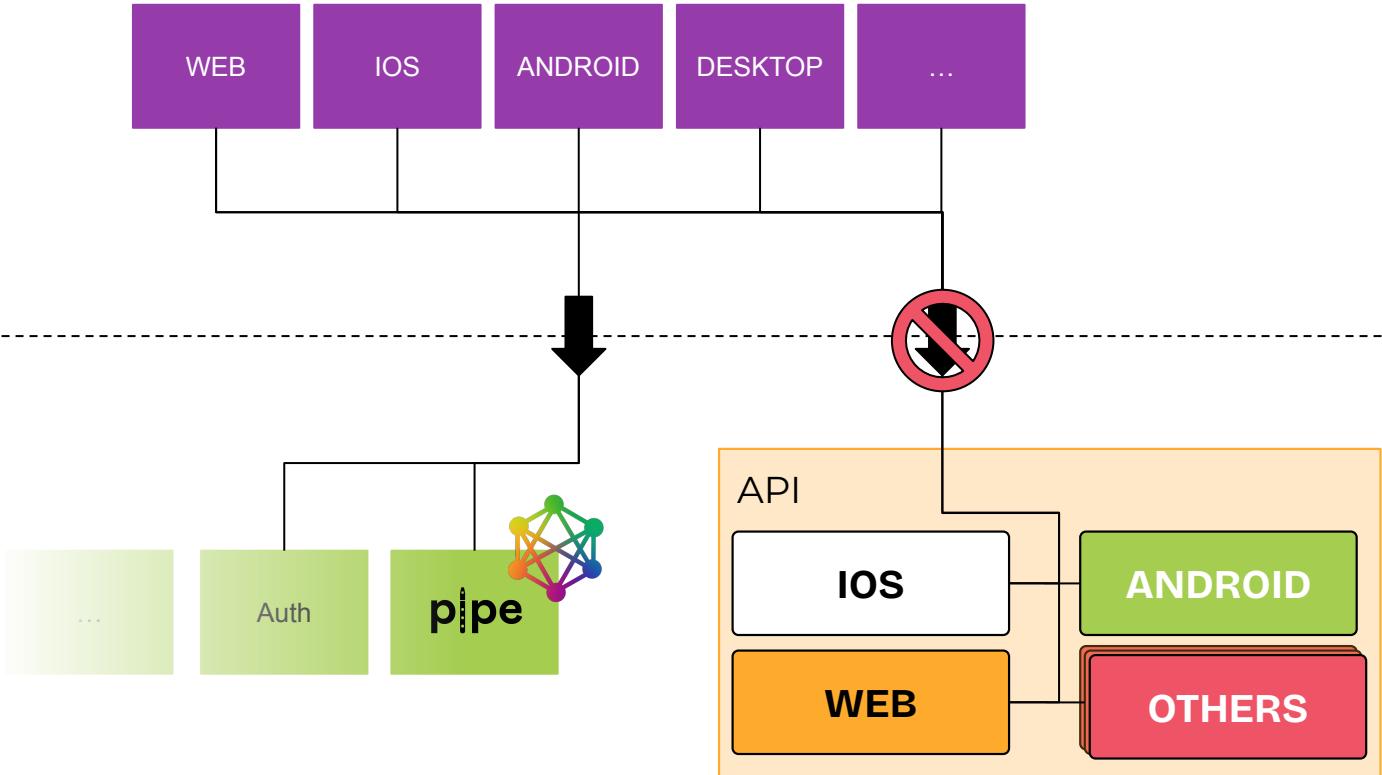
internal  
Client & Devices



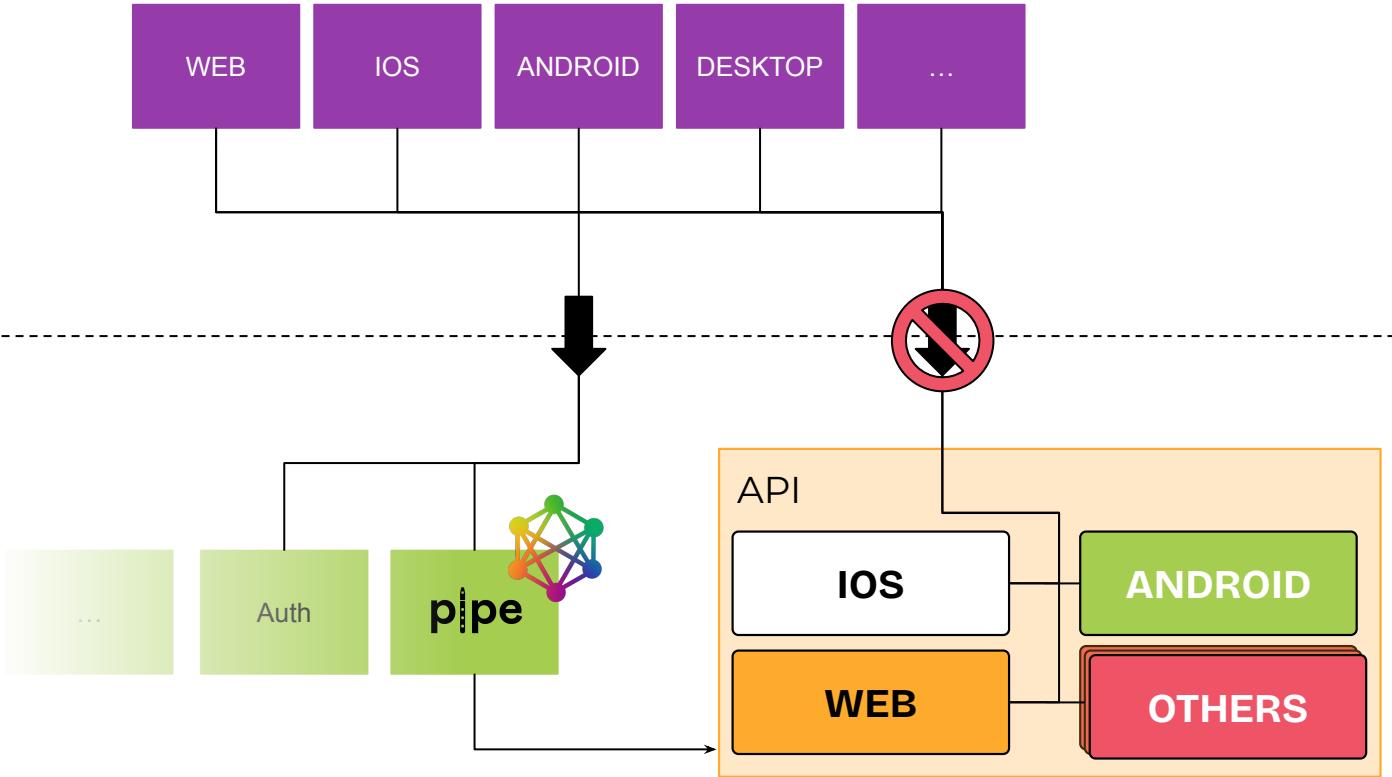
Client & Devices  
internal



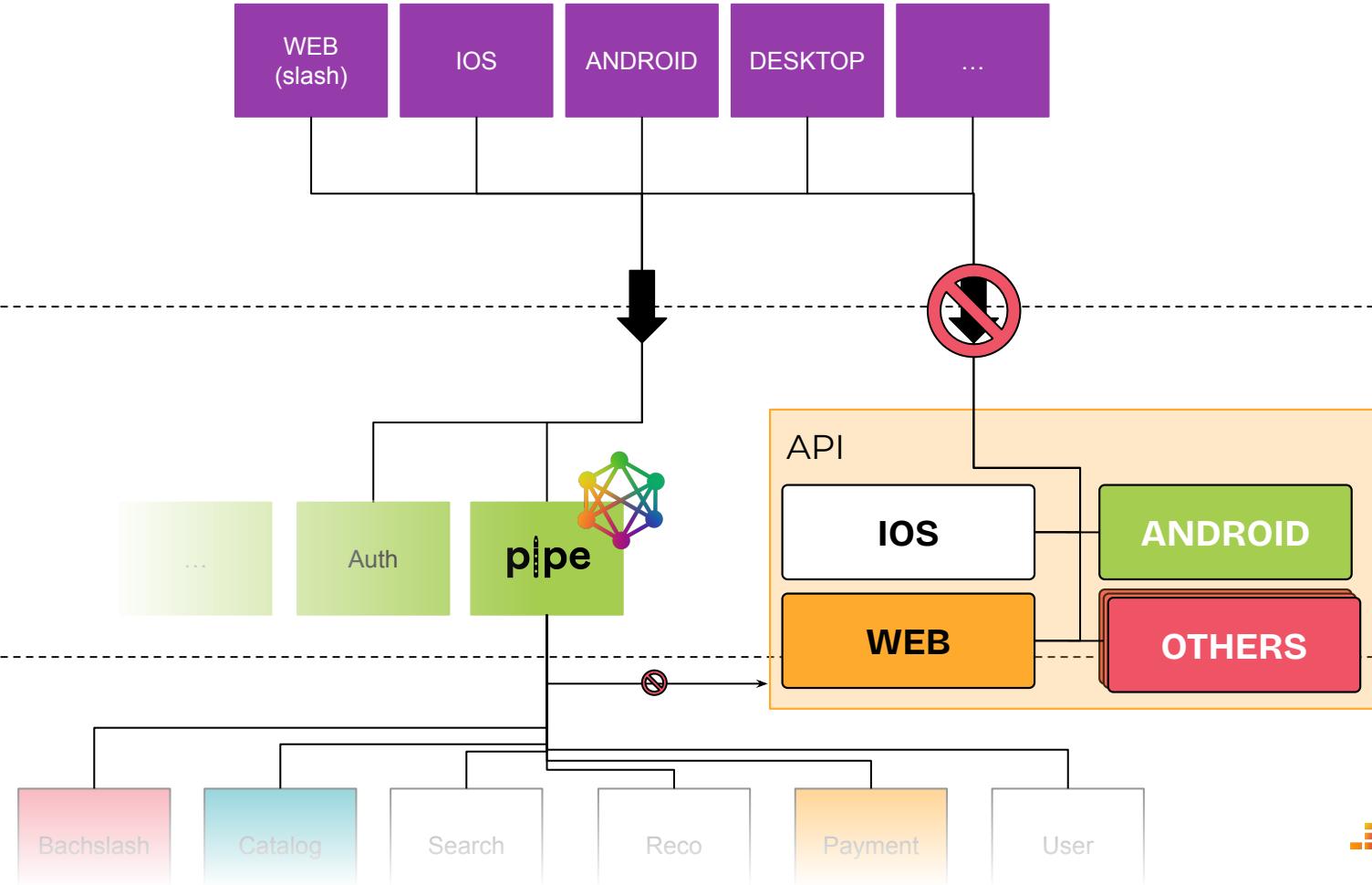
# Client & Devices



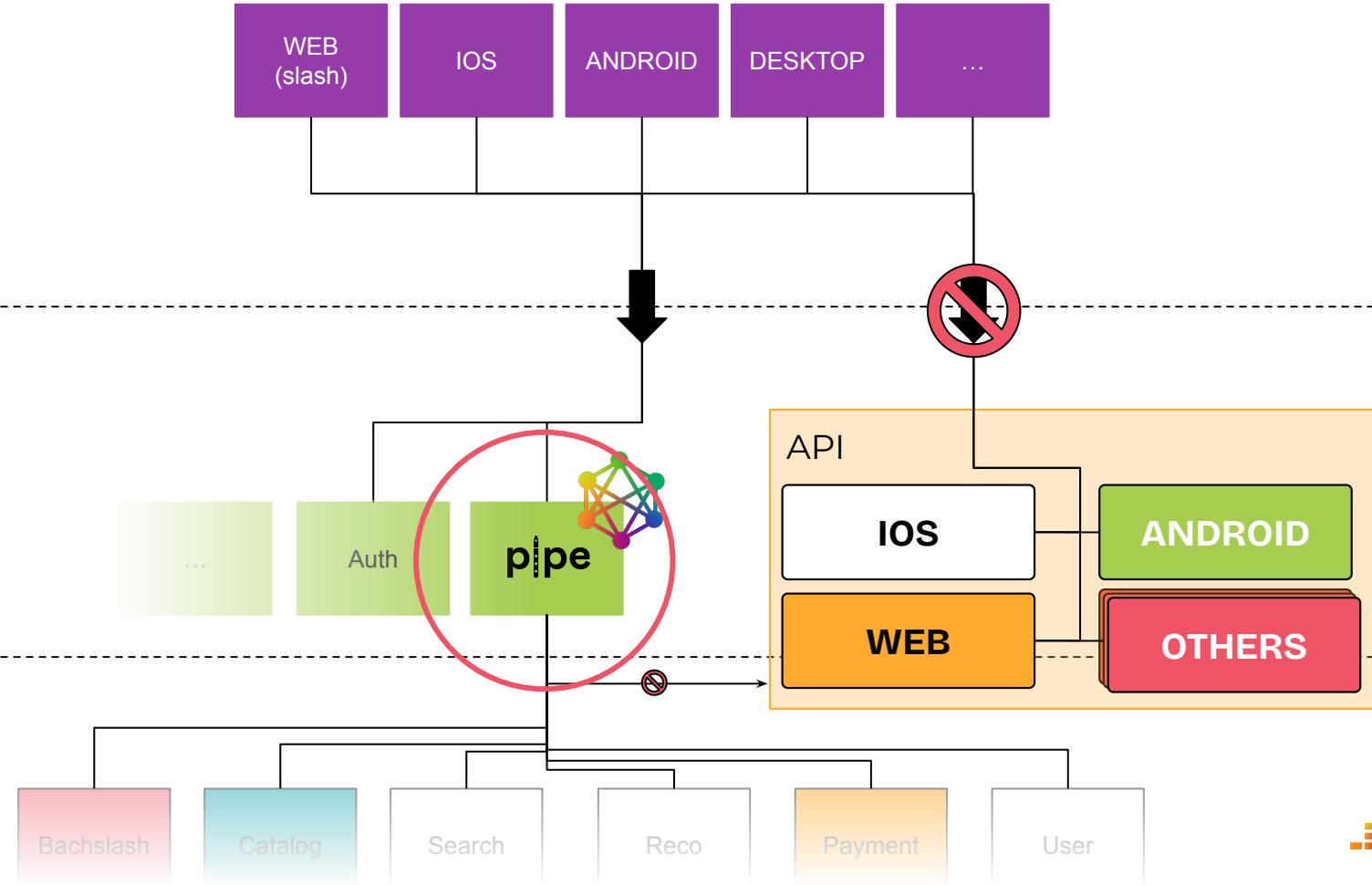
# Client & Devices



Client & Devices  
exposition  
internal



Client & Devices  
exposition  
internal

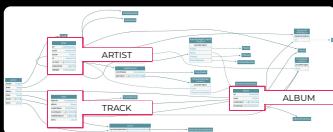
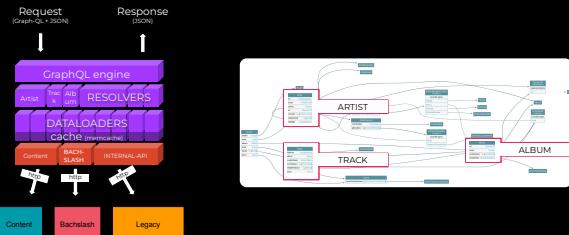




**Building**

# Multiple phases

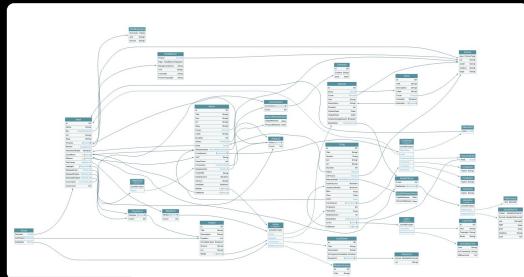
## Bootstrap



# Multiple phases

Bootstrap

Spread

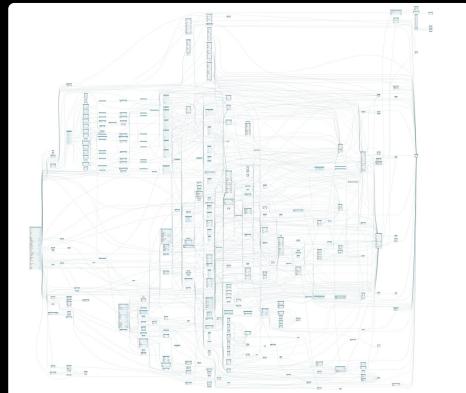


# Multiple phases

Bootstrap

Spread

Cruise



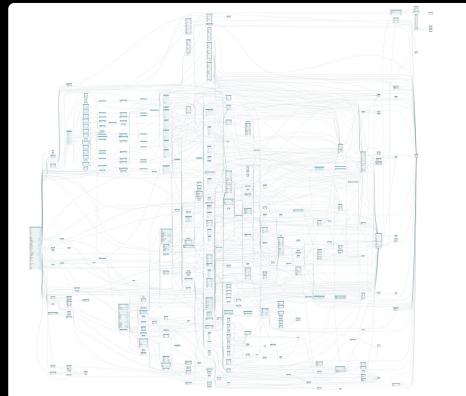
# Multiple phases

NOW

Bootstrap

Spread

Cruise

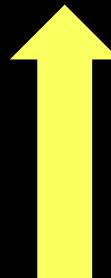


# Multiple phases

Bootstrap

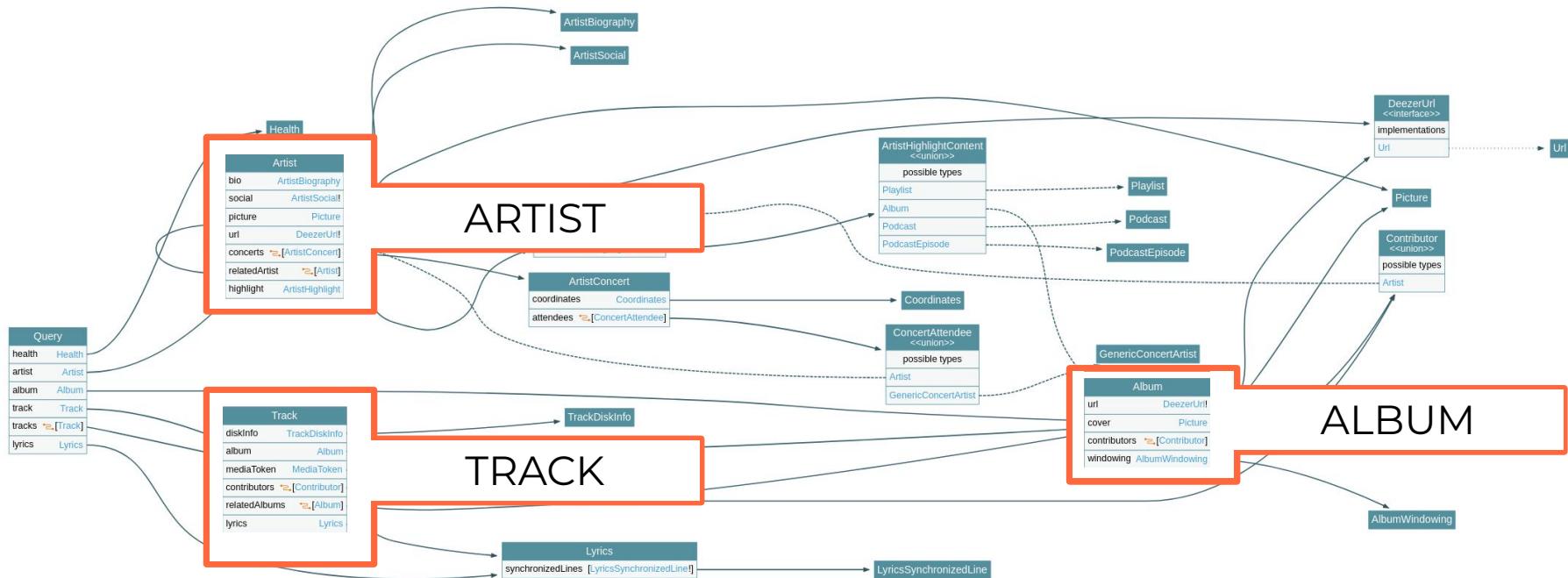
Spread

Cruise

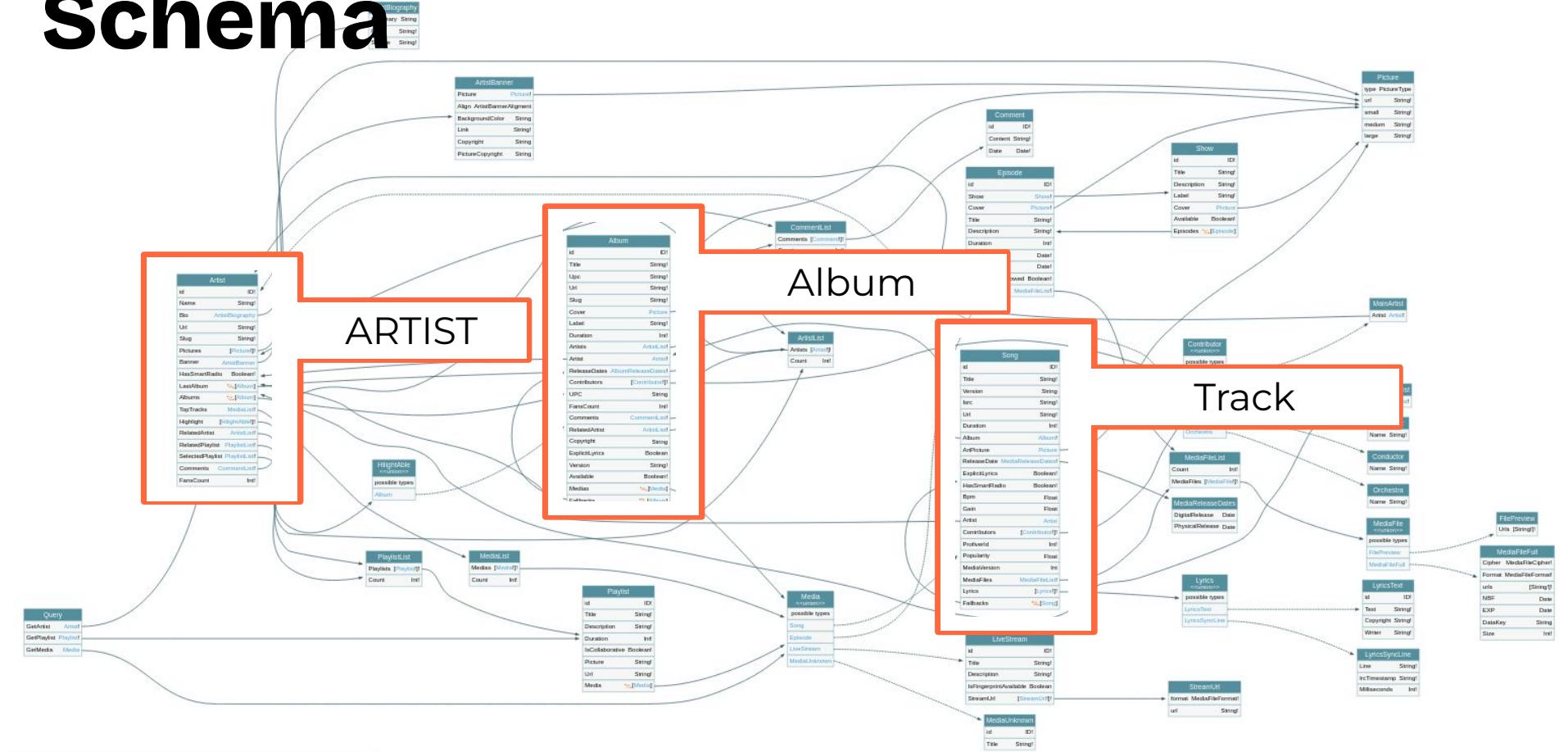


# Schema

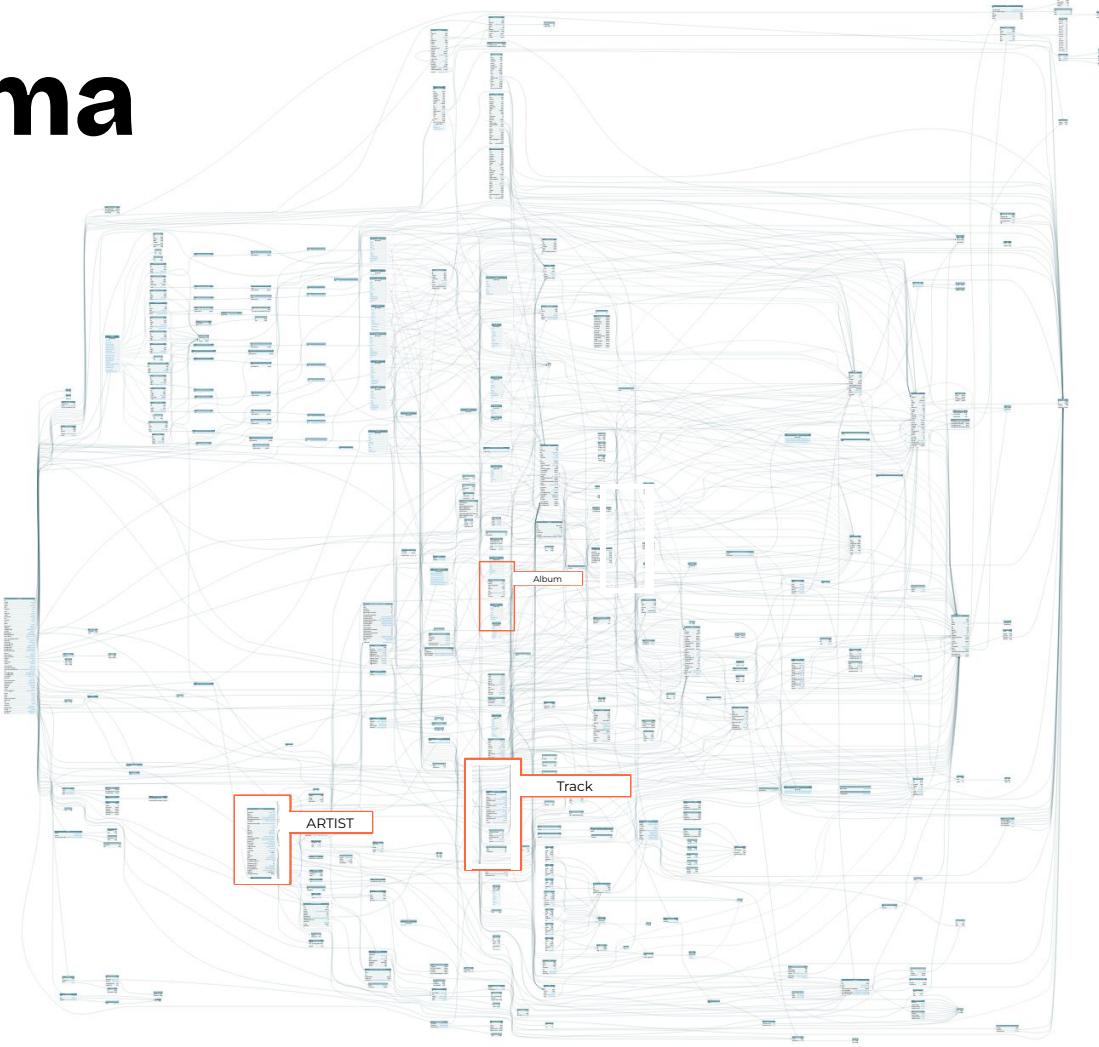
# Schema

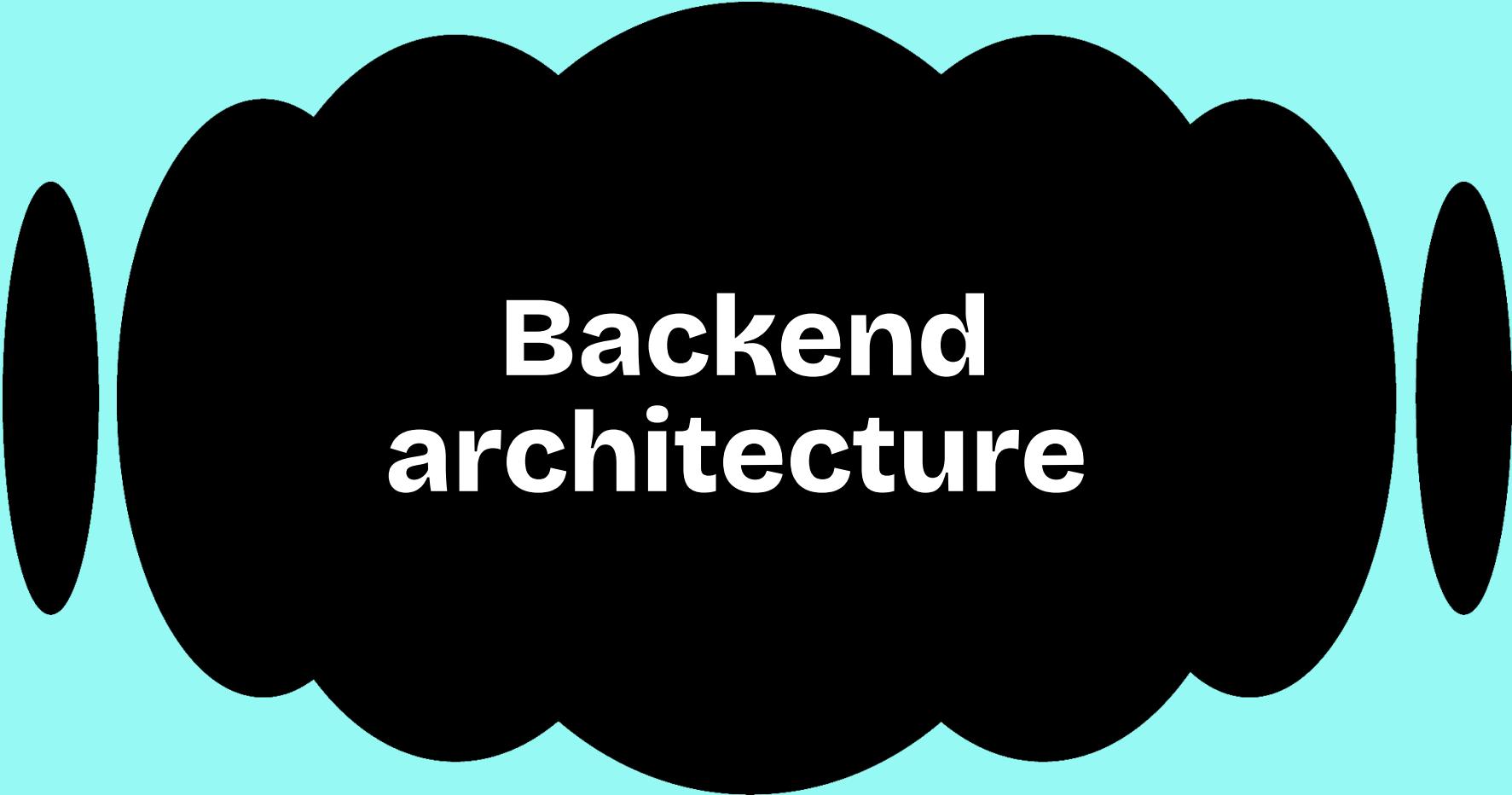


# Schema



# Schema

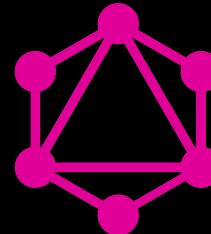
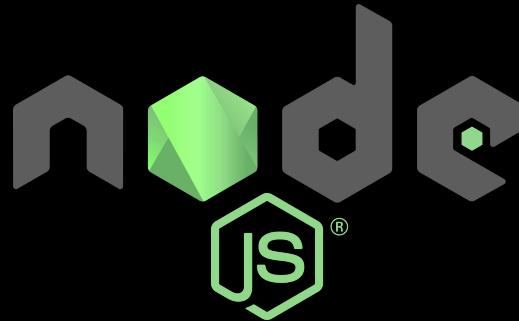




Backend  
architecture

# Backend architecture

- NodeJS 10 then 12 (we've upgraded since 😊)
- Typescript
- Express
- Standard GraphQL engine from Facebook
- Classic Resolver/Dataloader
- Gateway making HTTP calls internally
- (almost)No business logic



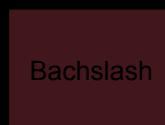
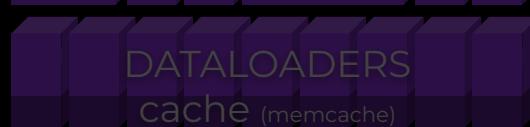
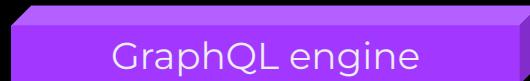
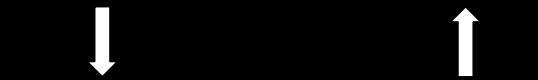
GraphQL

# Backend architecture

- NodeJS 10 then 12 (we've upgraded since 😊)
- Typescript
- Express
- Standard GraphQL engine from Facebook
- Classic Resolver/Dataloader
- Gateway making HTTP calls internally
- (almost)No business logic

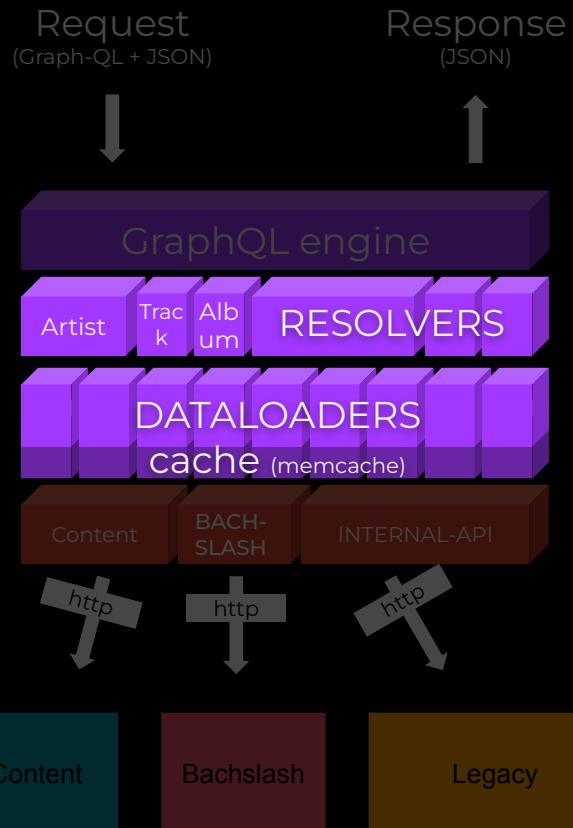
Request  
(Graph-QL + JSON)

Response  
(JSON)



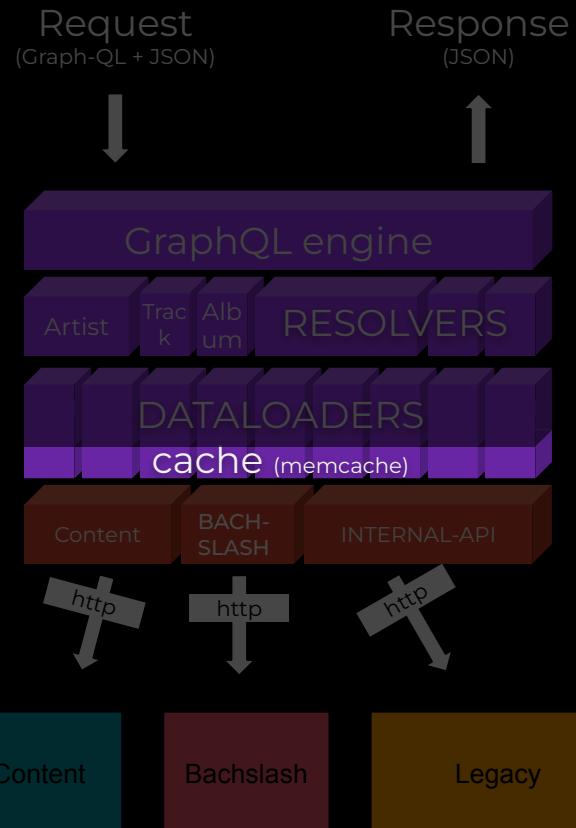
# Backend architecture

- NodeJS 10 then 12 (we've upgraded since 😊)
- Typescript
- Express
- Standard GraphQL engine from Facebook
- Classic Resolver/Dataloader
- Gateway making HTTP calls internally
- (almost)No business logic



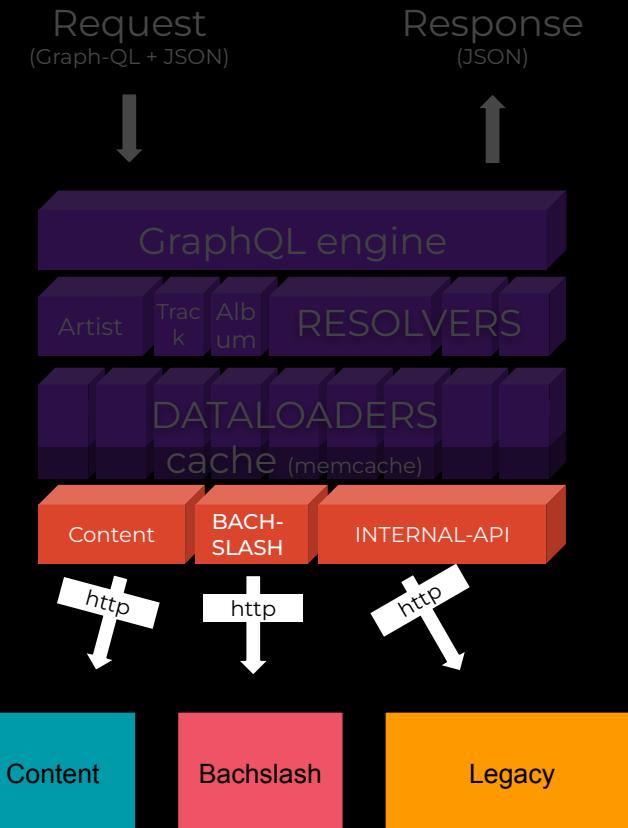
# Backend architecture

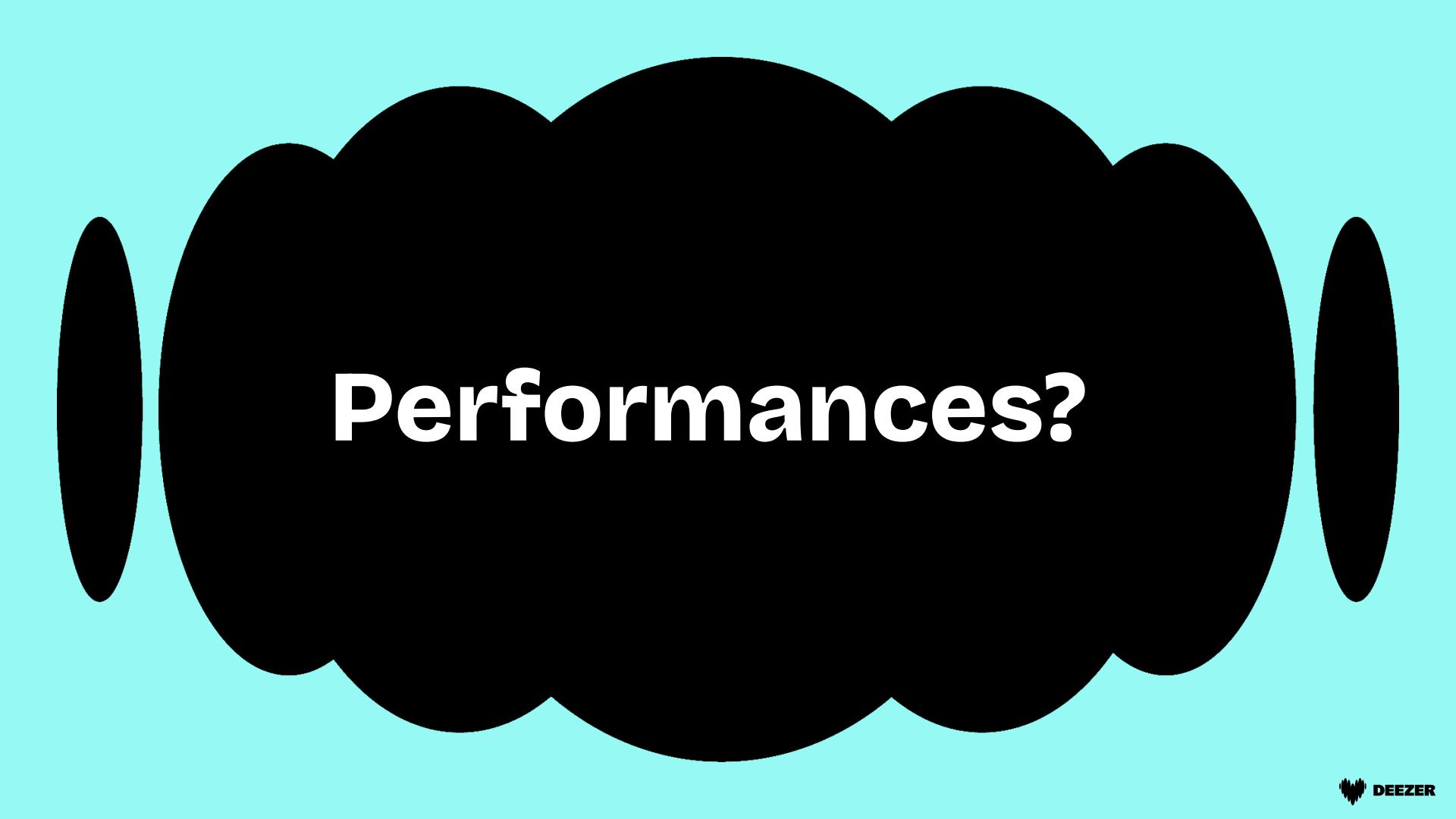
- NodeJS 10 then 12 (we've upgraded since 😊)
- Typescript
- Express
- Standard GraphQL engine from Facebook
- Classic Resolver/Dataloader
- Gateway making HTTP calls internally
- (almost)No business logic



# Backend architecture

- NodeJS 10 then 12 (we've upgraded since 😊)
- Typescript
- Express
- Standard GraphQL engine from Facebook
- Classic Resolver/Dataloader
- Gateway making HTTP calls internally
- (almost)No business logic





**Performances?**

# How fast was it?

To load an artist



To load an album

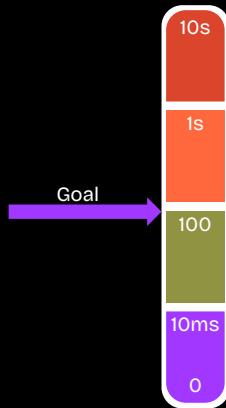


To load 50 tracks



# How fast was it?

To load an artist



To load an album

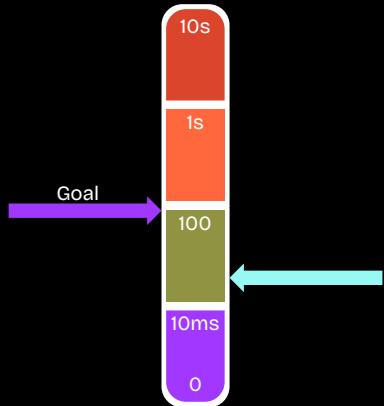


To load 50 tracks

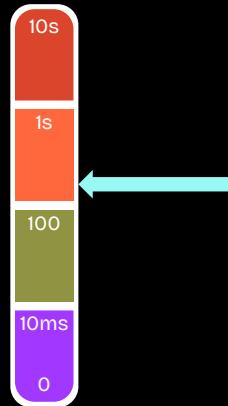


# How fast was it?

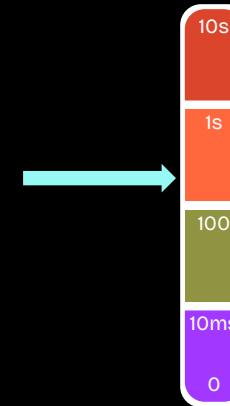
To load an artist



To load an album

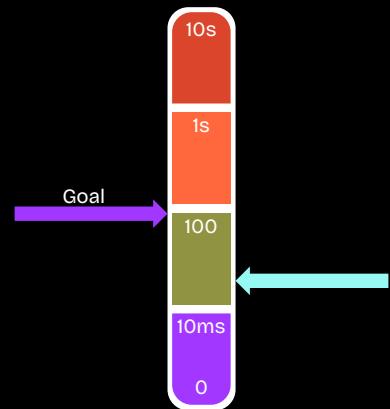


To load 50 tracks

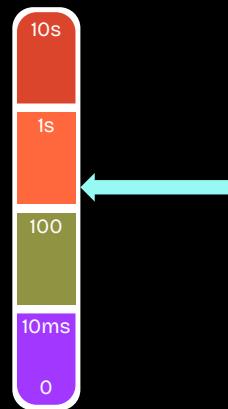


# How fast was it?

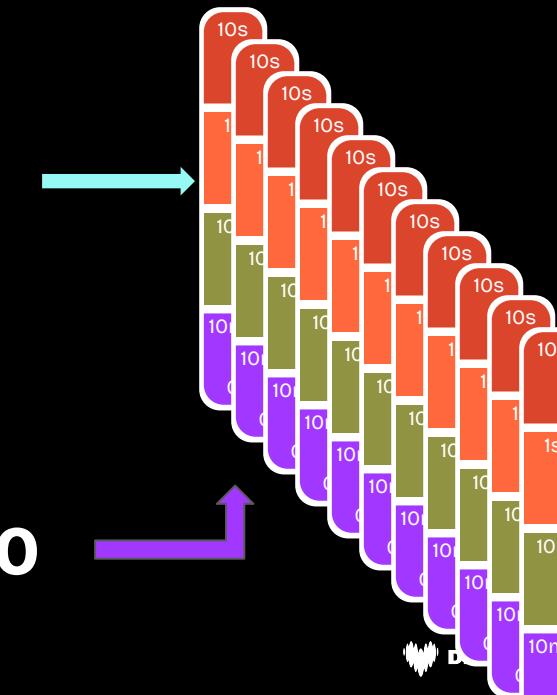
To load an artist



To load an album



To load 50 tracks



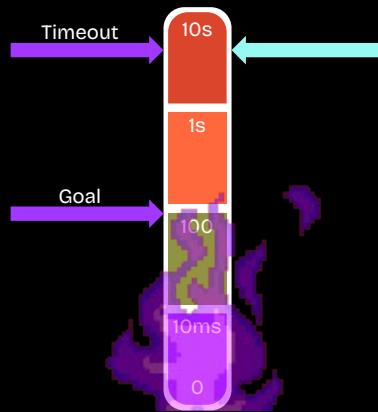
Loading a playlist of 1,000 tracks →

x20

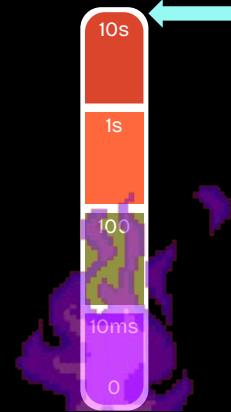
# How fast was it?



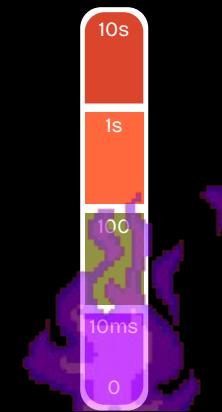
To load an artist



To load an album

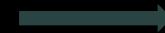


To load 50 tracks



Under (not so)serious load

# How fast was it ?

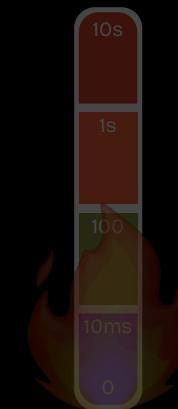


To load an Arti



The Slow Mo Guys

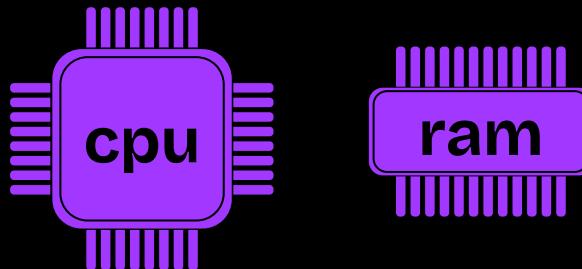
oad 50 tracks



Under (not so)serious load

# How hungry is it?

**How many servers** do I need to  
ask the Infra team to buy or re-use?

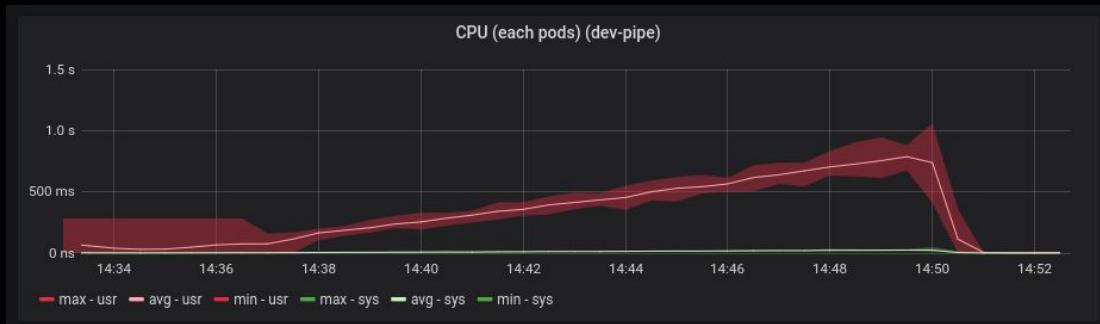


# How hungry is it?

CPU



## What resources to monitor?

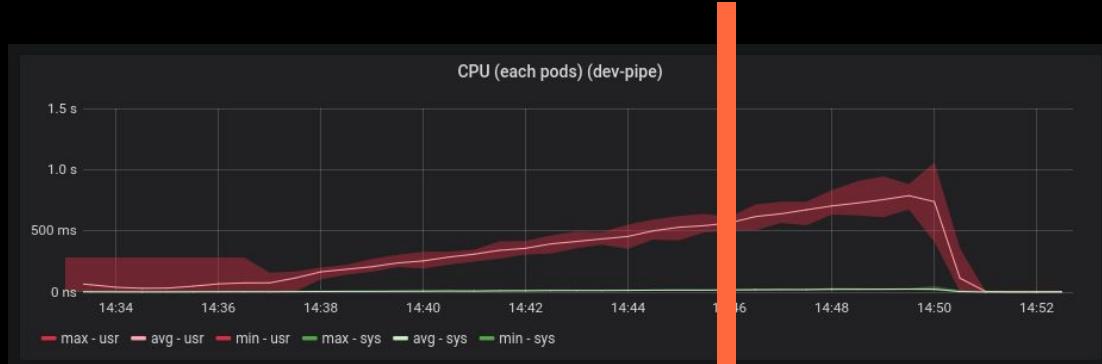


# How hungry is it?

CPU



What resources to monitor?



response time becomes bad

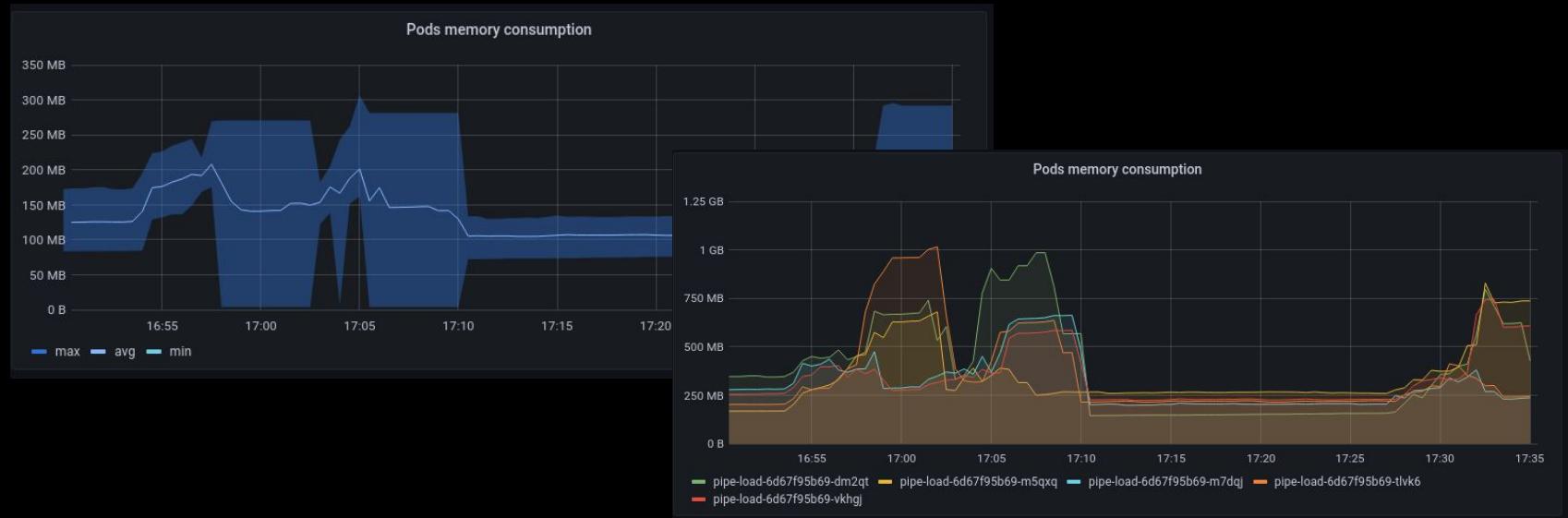
# How hungry is it?

What resources to monitor?

CPU



RAM



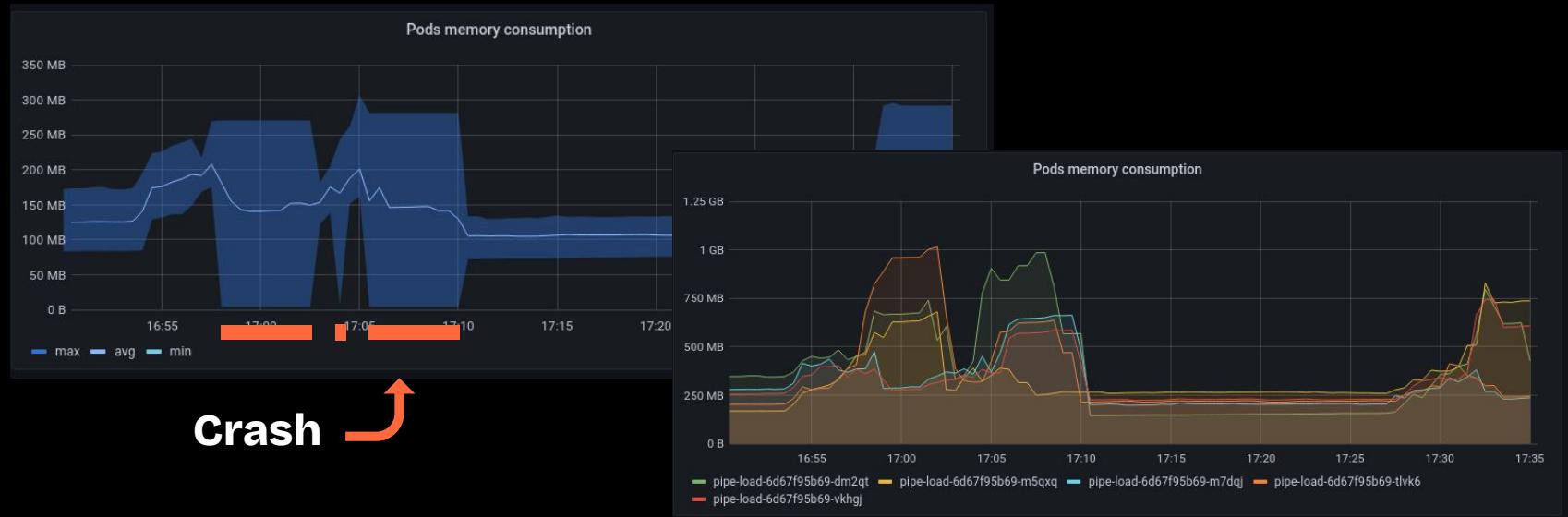
# How hungry is it?

What resources to monitor?

CPU

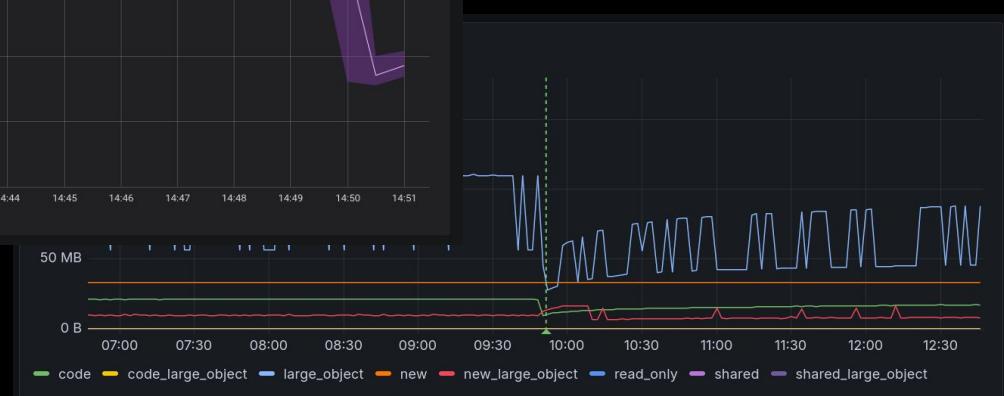
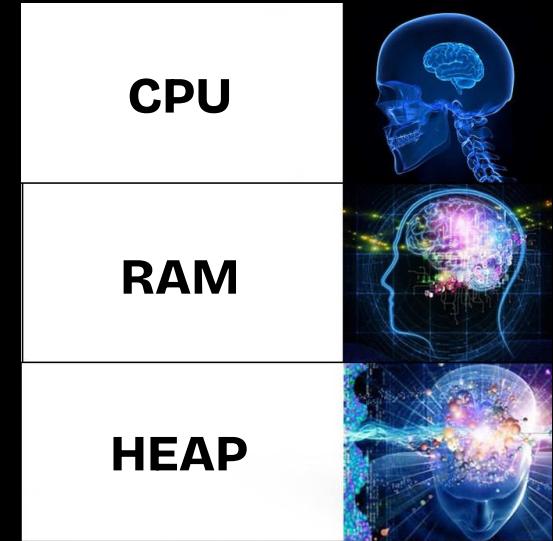
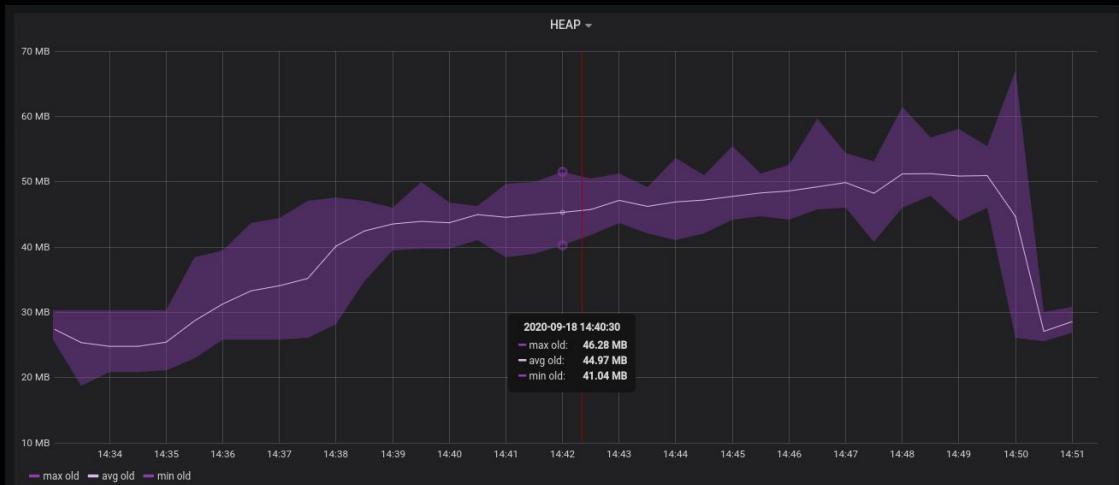


RAM



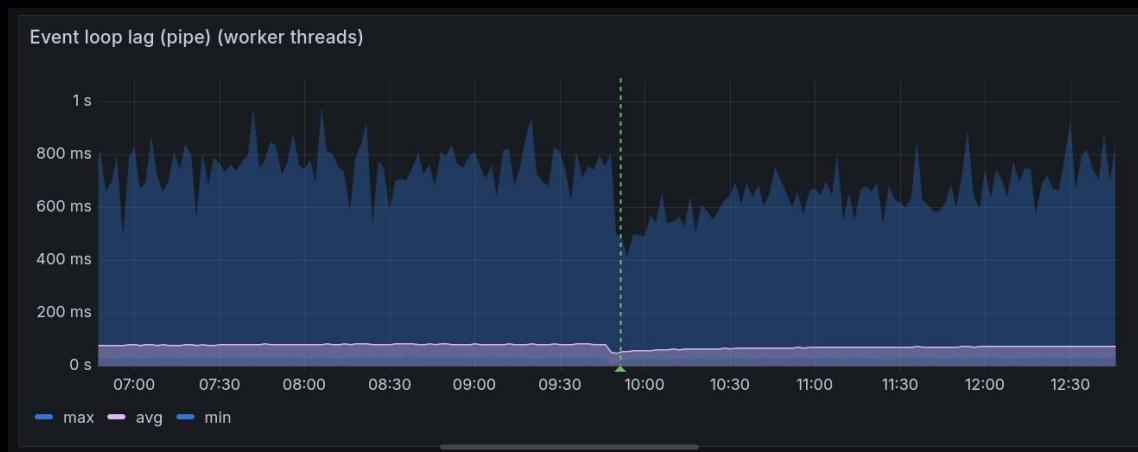
# How hungry is it?

What resources to monitor?



# How hungry is it?

What resources to monitor?



CPU



RAM



HEAP

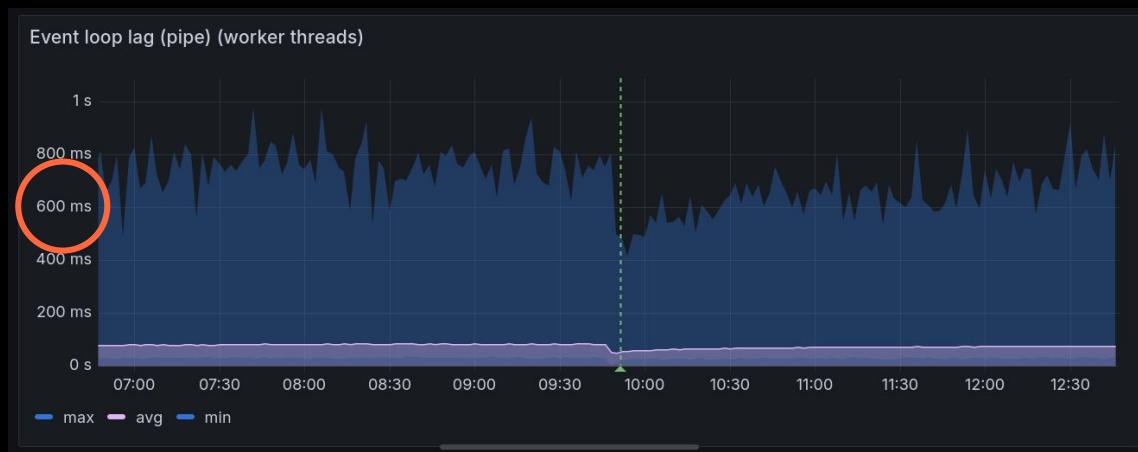


Event  
Loop lag



# How hungry is it?

What resources to monitor?



CPU



RAM



HEAP

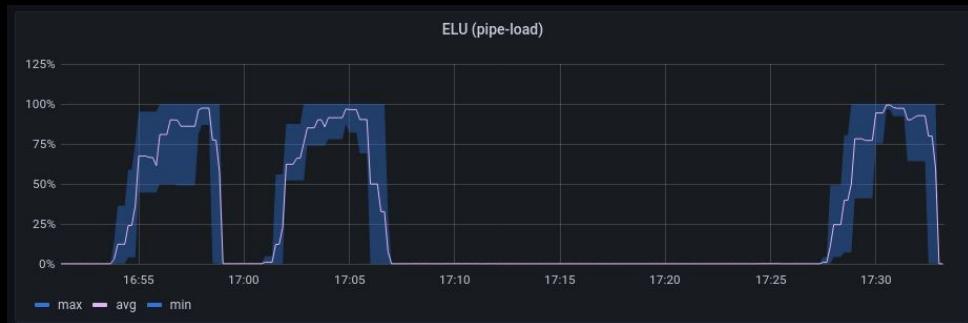


Event  
Loop lag



# How hungry is it?

What resources to monitor?



CPU



RAM



HEAP



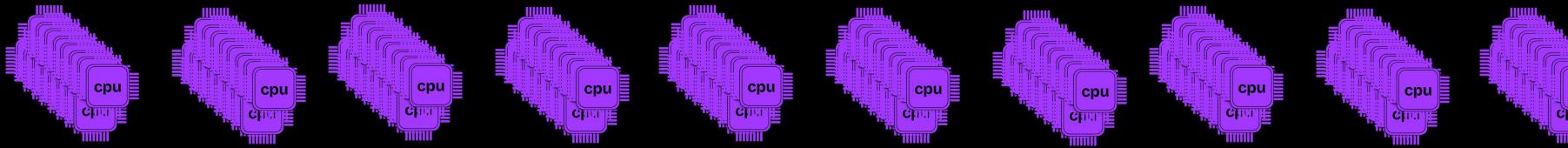
Event  
Loop lag



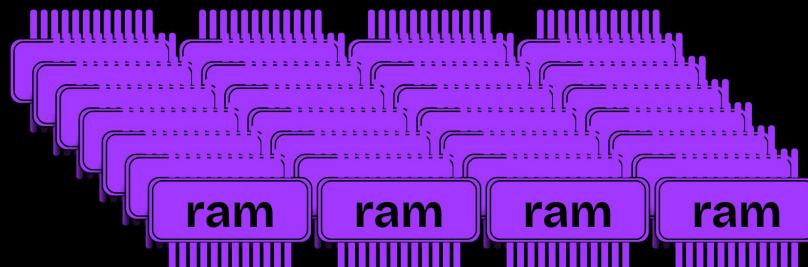
Event  
Loop  
Utilization



# Initial performances insight



More CPUs than the entire fleet of existing servers  
&  
Half the RAM of existing servers

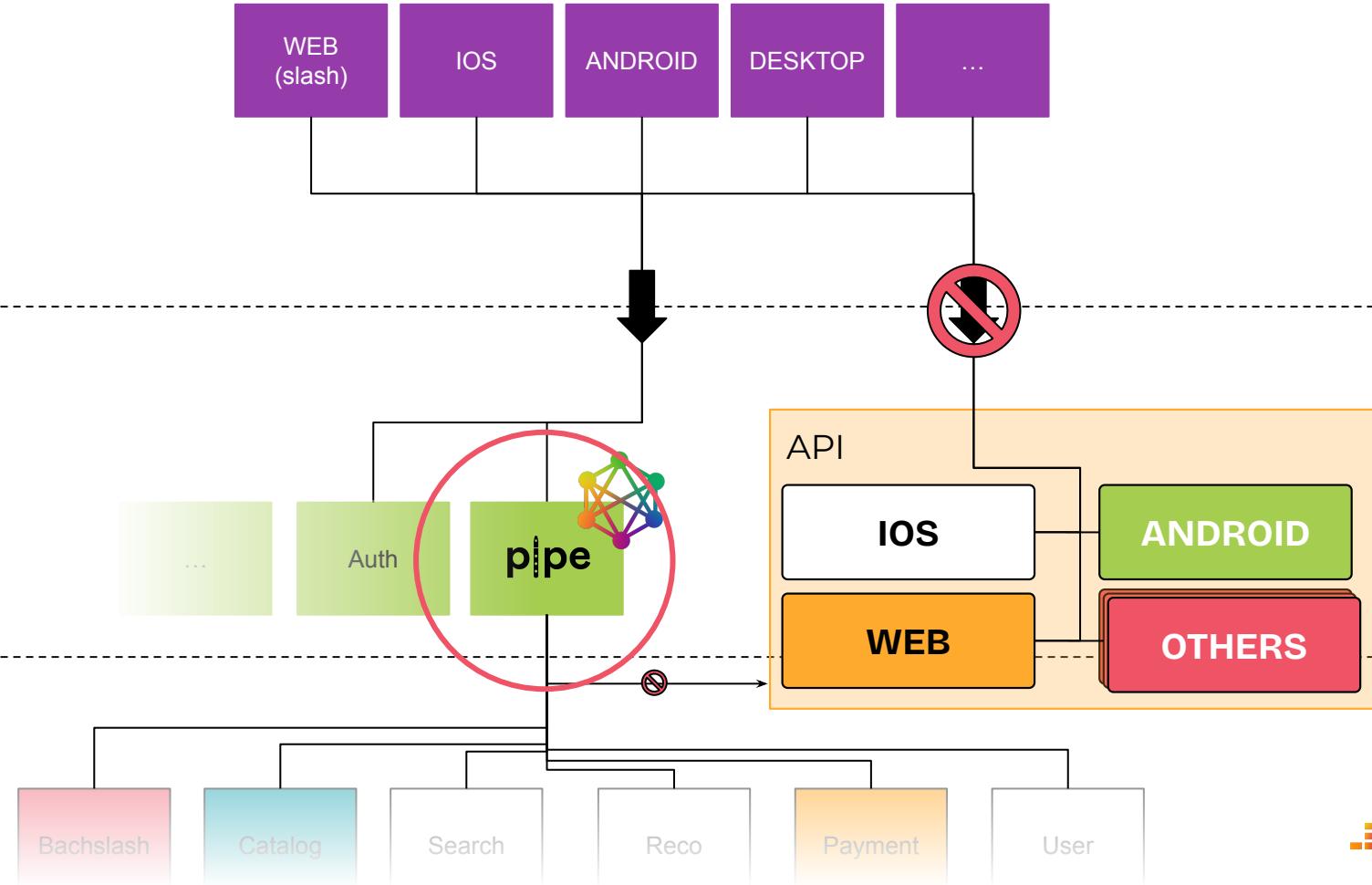


# Initial performances insight



The Slow Mo Guys

Client & Devices  
exposition  
internal

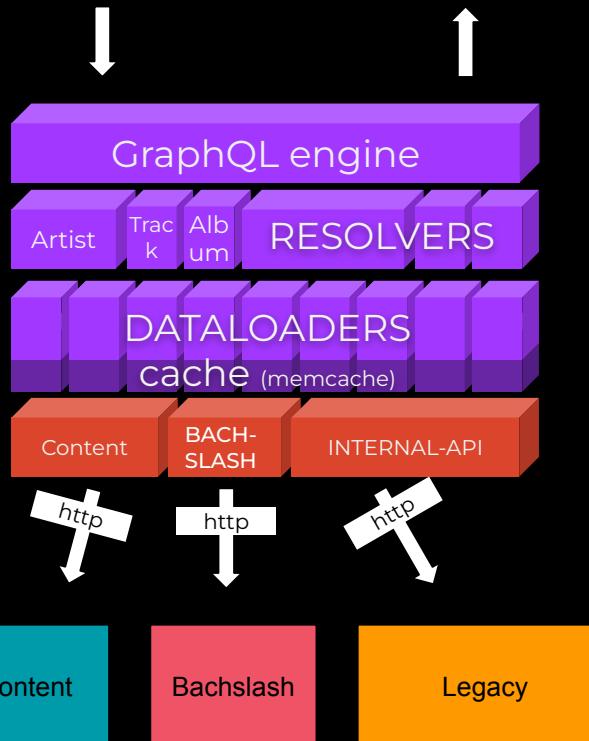


# Optimizing

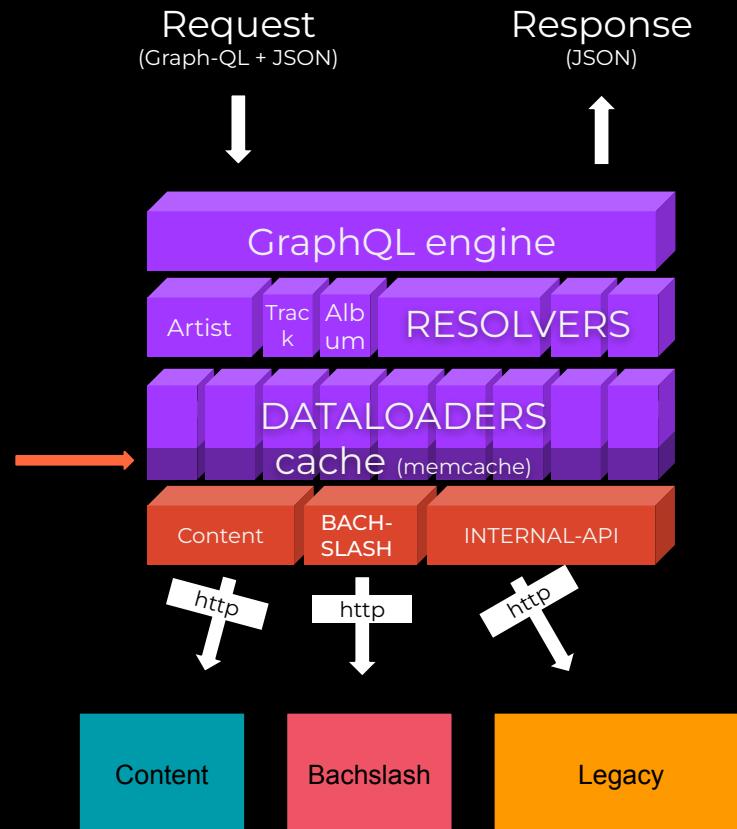
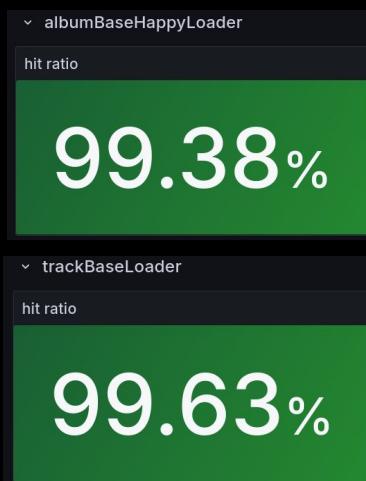
# Where to optimize?

Request  
(Graph-QL + JSON)

Response  
(JSON)

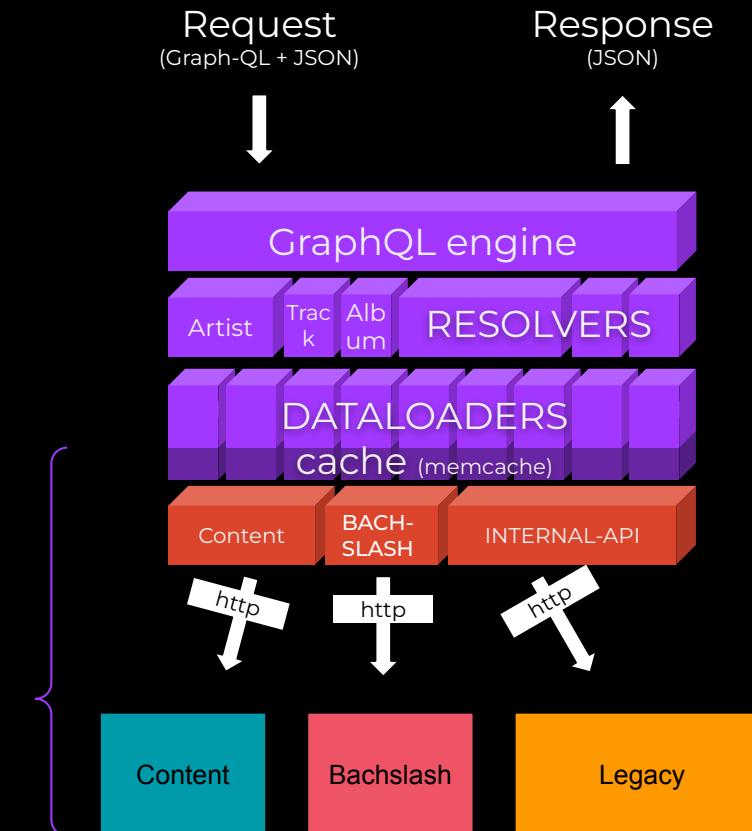


# Where to optimize?



# Where to optimize?

Cache allow us to skip the bottom part most of the time



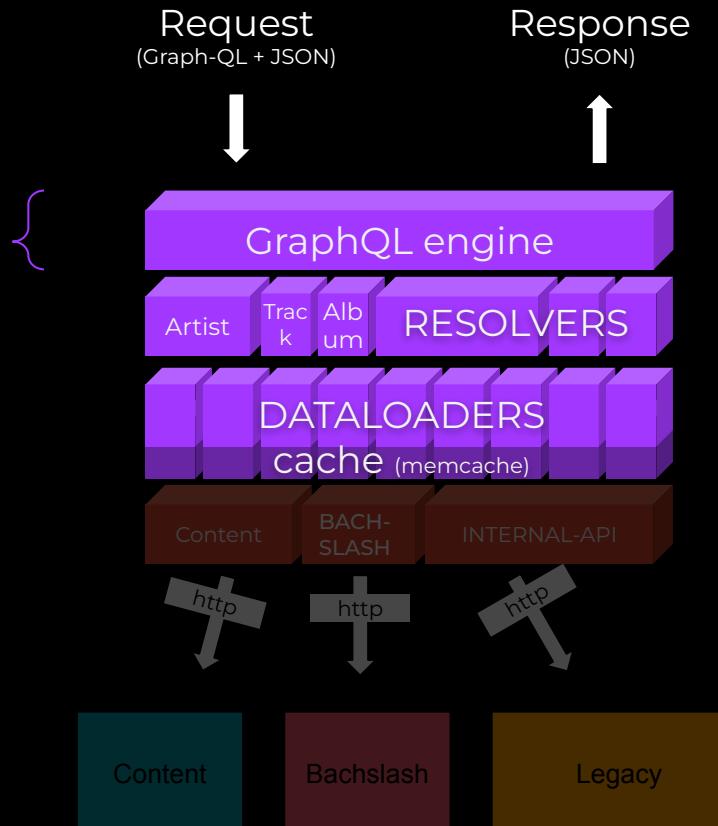
# Where to optimize?

GraphQL engine is not ours

- GraphQL JIT was promising and we tested it.

Turns out it is not worth it for us.

More info here: [Deezer.io - GraphQL JIT, is it worth it?](#)



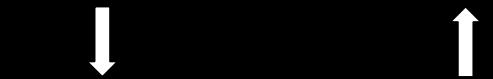
# Where to optimize?

Mainly

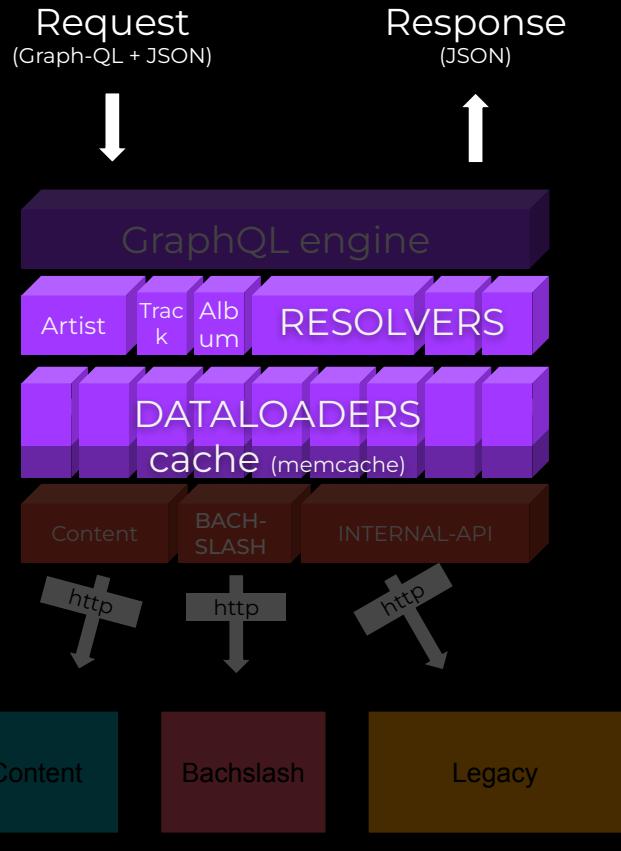
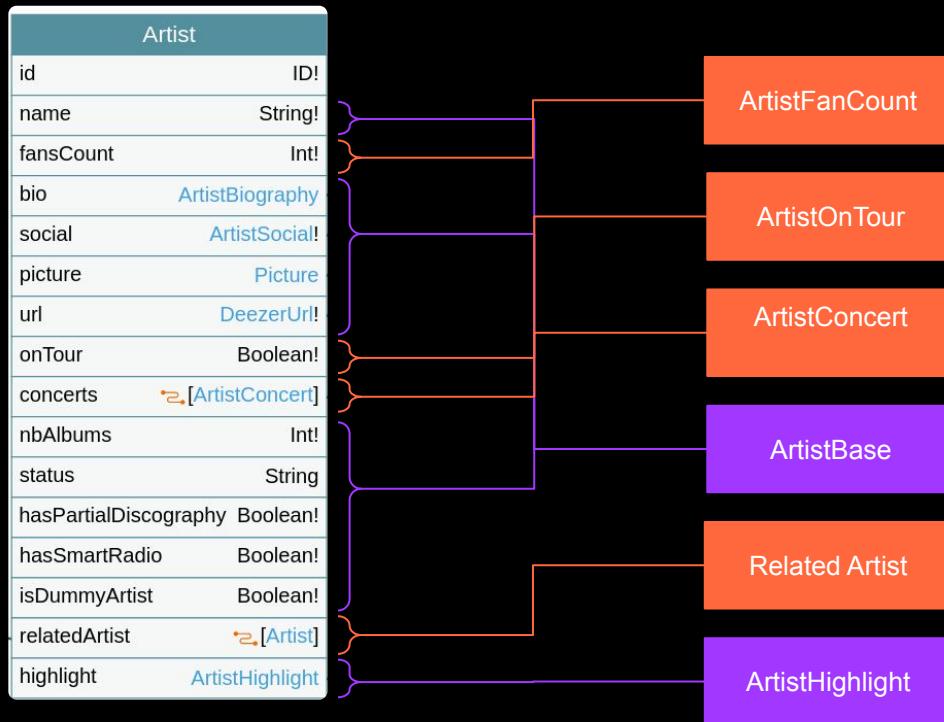
- Resolvers
- Dataloaders

Request  
(Graph-QL + JSON)

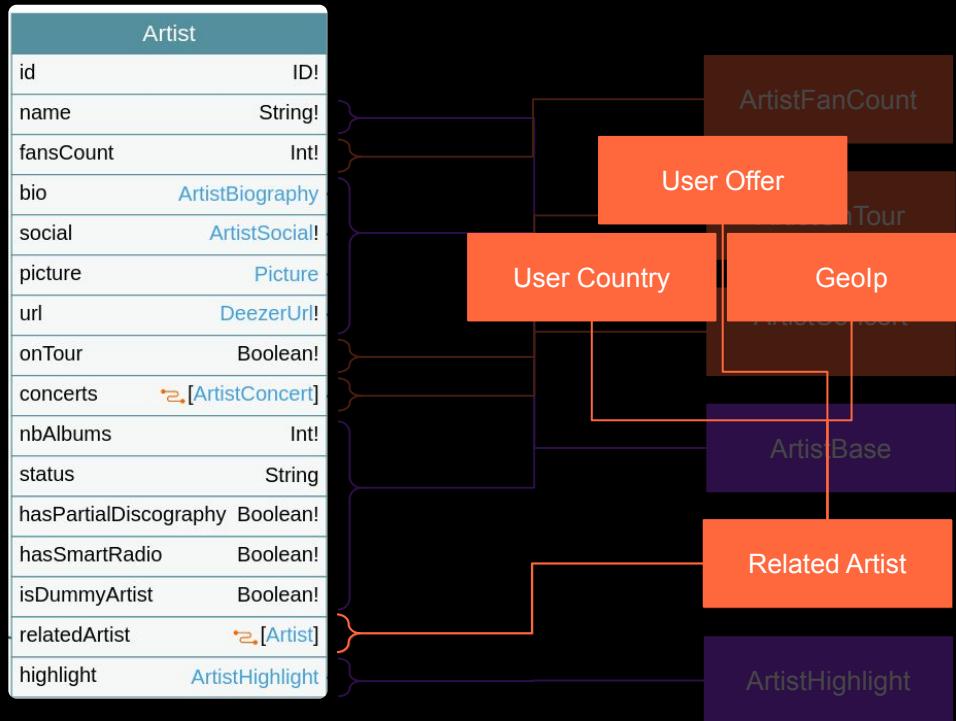
Response  
(JSON)



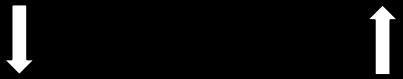
# Where to optimize?



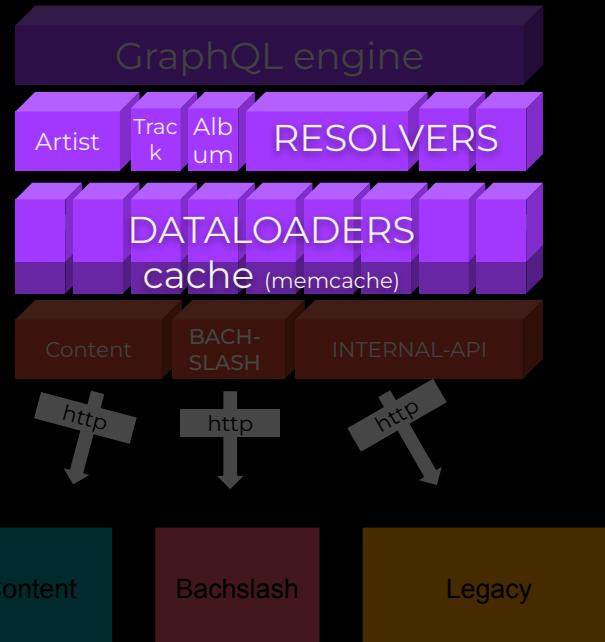
# Where to optimize?



Request  
(Graph-QL + JSON)



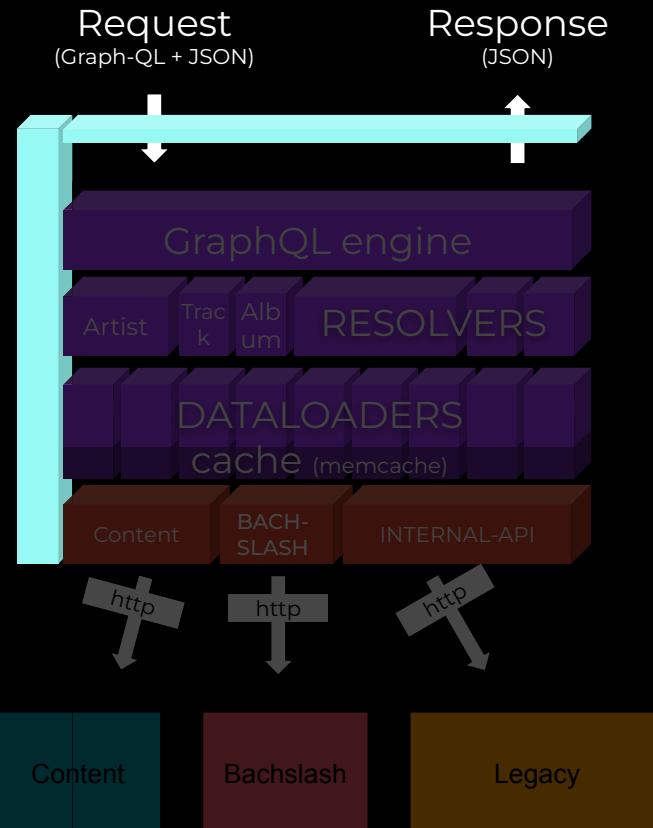
Response  
(JSON)



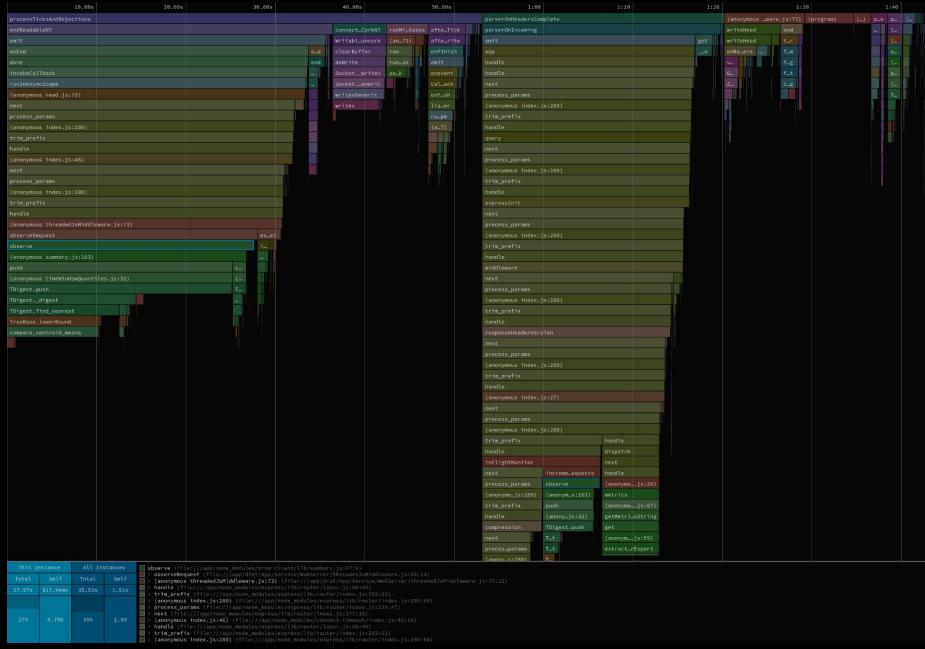
# Where to optimize?

All the utility code

- Monitoring
- Express middleware
- GraphQL middleware
- GraphQL Shield

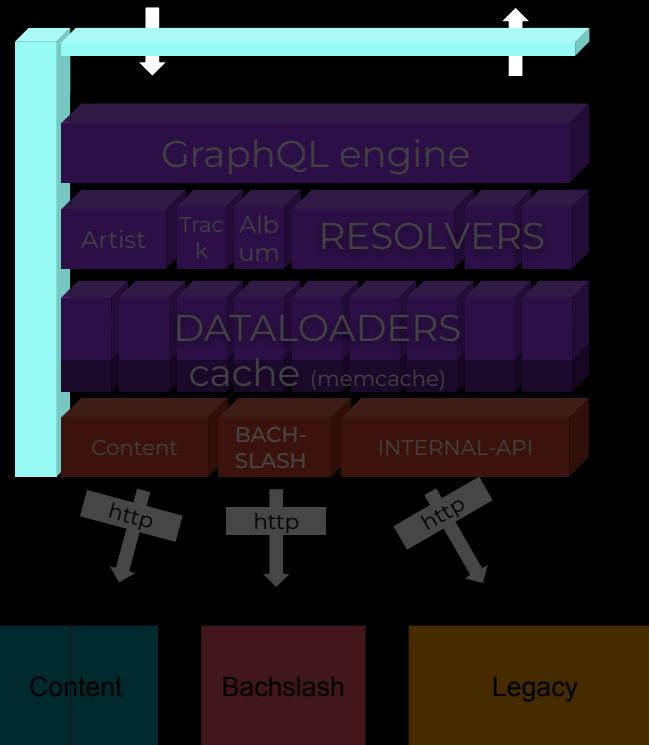


# Where to optimize?



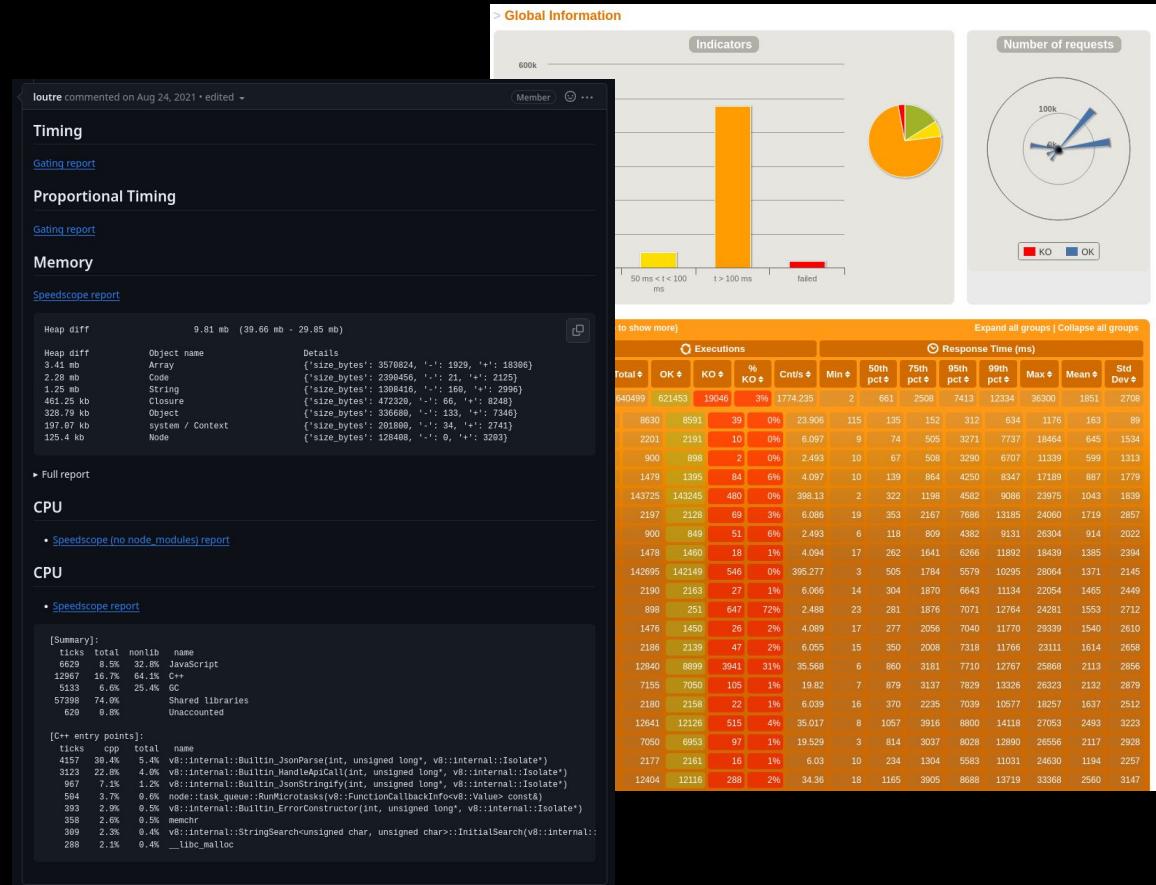
Request  
(Graph-QL + JSON)

Response  
(JSON)



# Where to optimize?

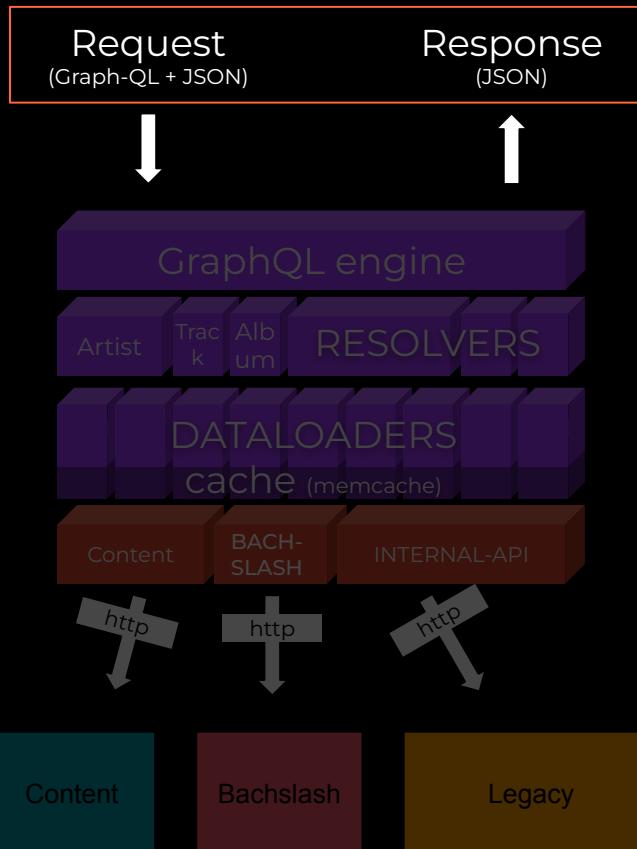
- Design realistic tests scenarios
- Automate them on each Pull Request
- Observe improvement PRs after PRs



# Where to optimize?

## Batching

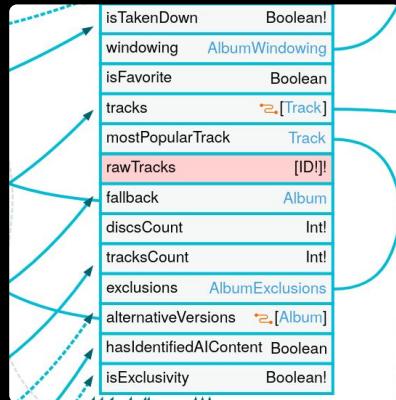
- Smaller sequential batches ( $500 \rightarrow 100$  or  $50$ )
- Better caching on client side
- Single entity calls for big ones (Artist/Album)



# Where to optimize?

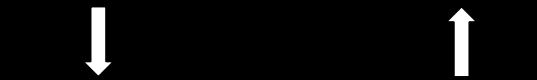
## Pagination

- Raw Tracks → Fast but no filtering
- Parallel query of smaller batches (50)



Request  
(Graph-QL + JSON)

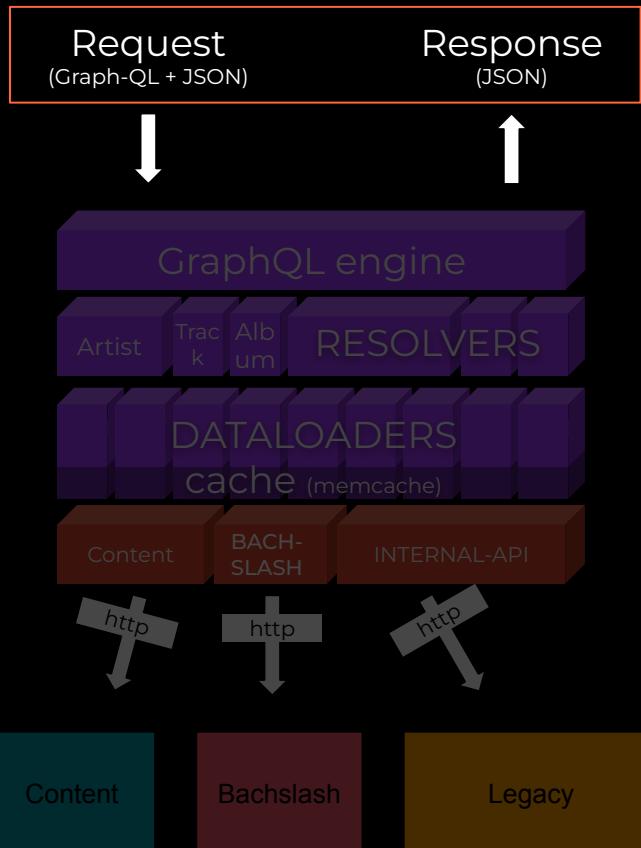
Response  
(JSON)



# Where to optimize?

## Smaller queries

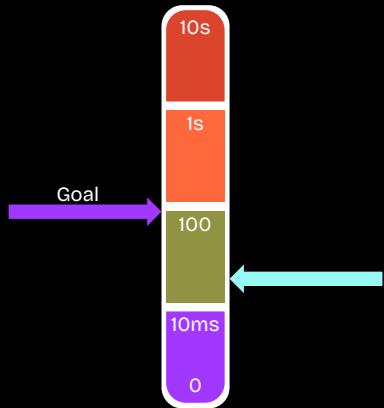
- Artist / Album / Tracks are fetch separately
- Client cache help reduce number of request
- More small requests spread load more easily



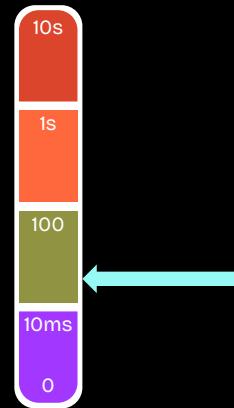
# Results

# How fast did it become?

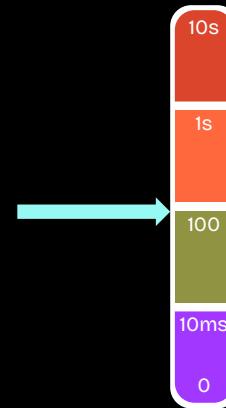
To load an artist



To load an album

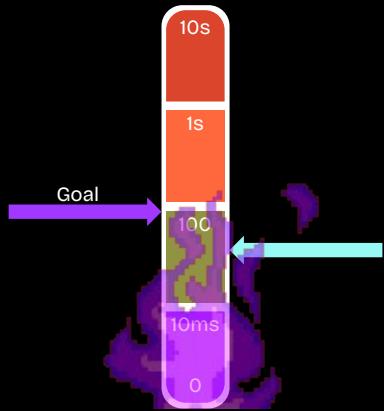


To load ~~50~~ 100 tracks

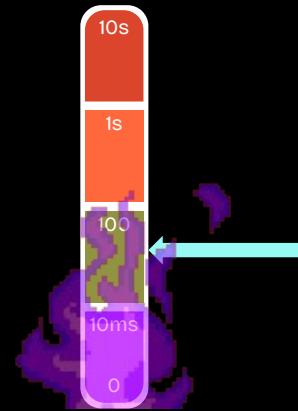


# How fast did it become?

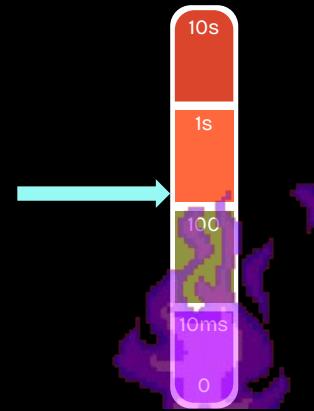
To load an artist



To load an album

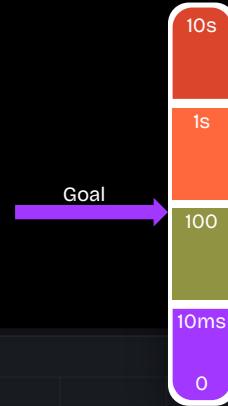


To load ~~50~~ 100 tracks



Under (real)serious load

# Everything together

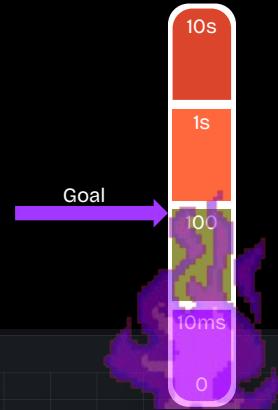


**93% of  
all requests  
under 100ms**



# Everything together

Under serious load

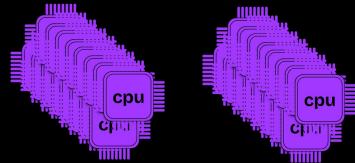


Goal

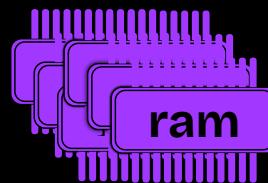
~70 to 80% of  
all requests  
under 100ms



# Everything together



only ~15% of the estimated  
resources





Takeaways

# Building a GraphQL API

## Collaboration

All stacks and clients need to collaborate during all the designs and test phases in order to avoid pitfalls

How the API is used matters

# Building a GraphQL API

## Collaboration

All stacks and clients need to collaborate during all the designs and test phases in order to avoid pitfalls

How the API is used matters

## Cache is paramount

You won't survive without Caching on both sides: Client & Backend

# Building a GraphQL API

## Collaboration

All stacks and clients need to collaborate during all the designs and test phases in order to avoid pitfalls

How the API is used matters

## Cache is paramount

You won't survive without Caching on both sides: Client & Backend

## NodeJS is fast but Graph can be big

The bigger the graph, the bigger the data in memory and JSON to parse and build.

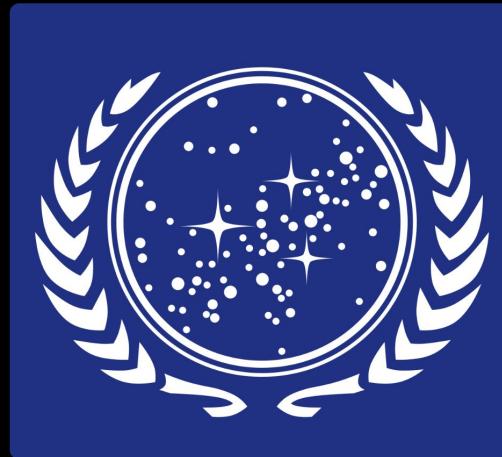
Parsing and producing JSON will be the last unavoidable bottleneck



**what's next?**

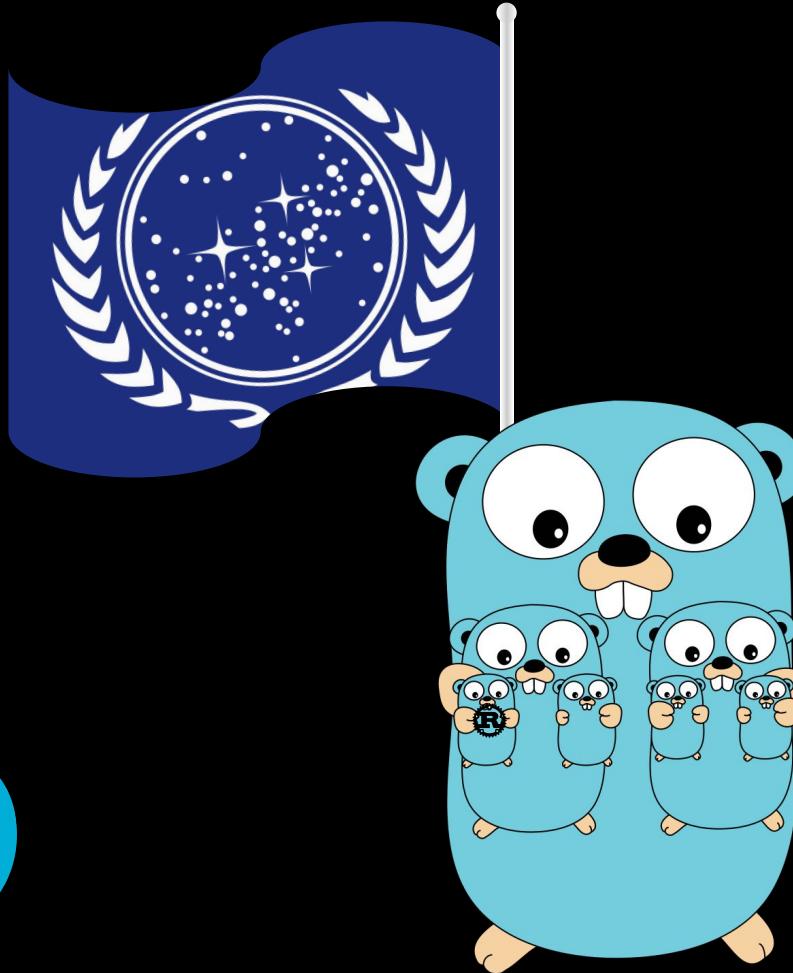
# Future plans

- **GraphQL Federation?**



# Future plans

- GraphQL Federation?
- Rewrite everything in GO? 😊



Q & A

# Thank you

