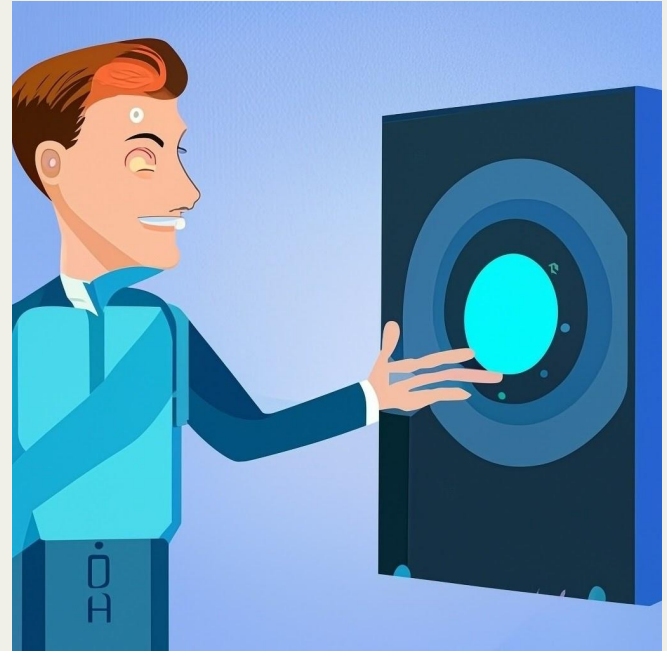# Multi-Tenancy!?

## Why should we care?

Imagine an app that lets you upload all the files on your hard drive and chat with it.
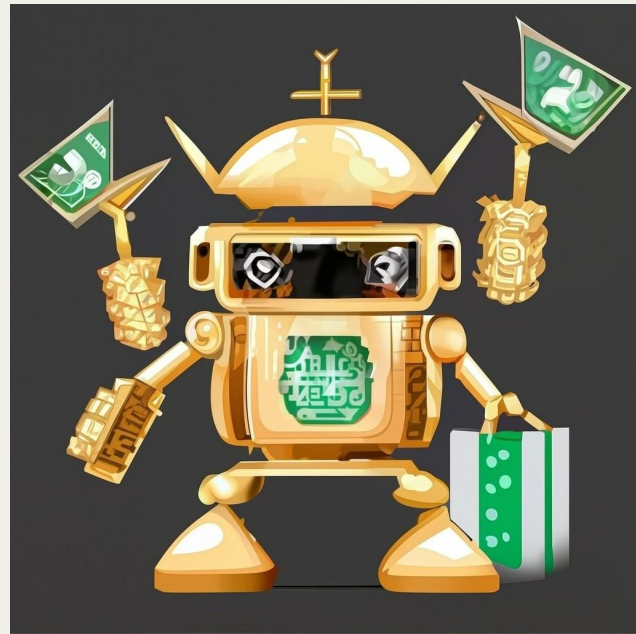
You wouldn't want someone else to chat with *your* files.

As the vendor of this chatbot app:

**Why should you pay for all your potential users if only 5% are active right now?**

**Previous multi-tenancy solutions made you pay for all, not for active tenants.**

# Prioritization is hard.

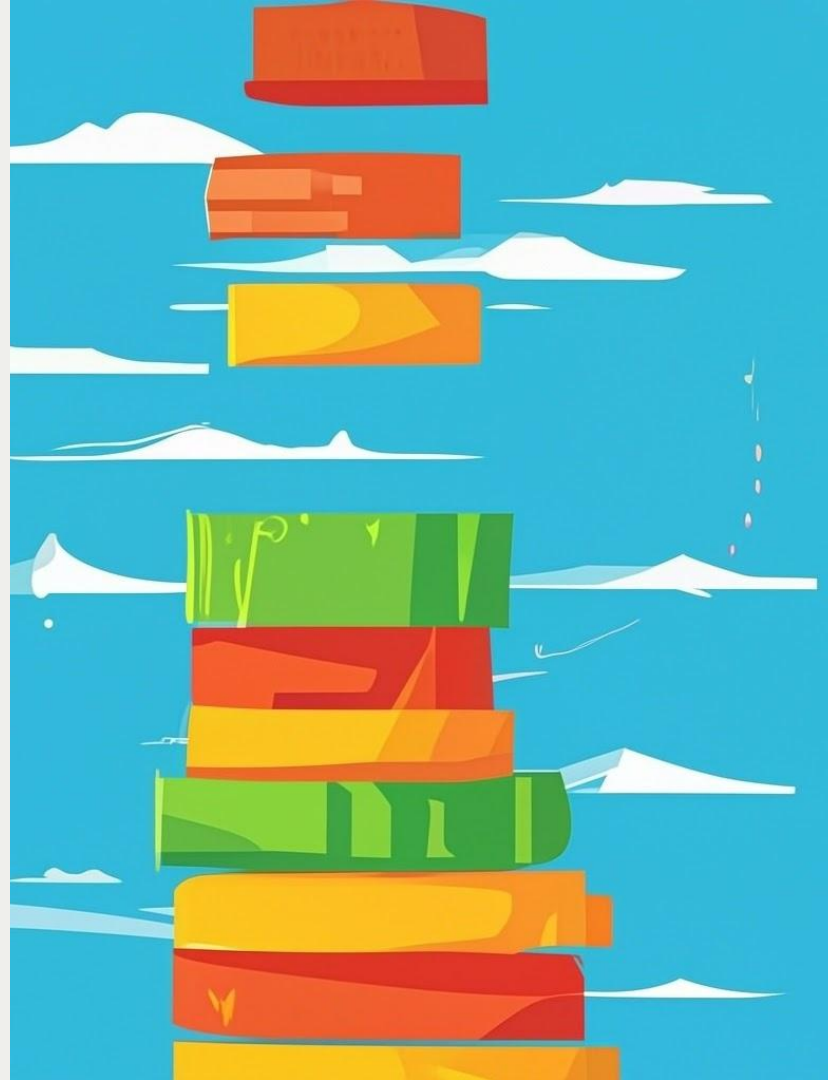## But sometimes it's obvious what you need to do.

**Commercial Customers:**

- "We need **strong isolation** between tenants for security and compliance (GDPR, etc.)"

- "We need to **reduce cost**."

## Engineers:

- "We need to add **hundreds of thousands of tenants** and there is no solution that can do this right now."

## Open–Source Community:

"I'm using workarounds with Weaviate and other vector DBs and they have limitations.

**Can you address the root cause?"**

# Why existing solutions failed

# Attempt Number 1

## The Filter

```
SELECT AWESOMENESS
FROM VECTORDB
WHERE TENANT_ID=12345
```

# Why is single index with filters not ideal?

**No isolation**

**Difficult to scale dynamically**

**Inefficient**

**Very expensive off-boarding**
(tenant deletion)

**Attempt Number 2**

# Create one collection per tenant

(The recommended way in Weaviate in the past)

👍

# Solves some of the previous shortcomings

Strict isolation

No filter necessary

Easy off-boarding

Good scalability

# Solves some of the previous shortcomings

Strict isolation

No filter necessary

Easy off-boarding

(In theory) Good scalability

# But isn't the right solution either

**Does not scale to millions of tenants**

👍 data scheduling scales well

👎 schema does not

**Lots of duplication**

**Terrible MTTR**

# How can we make this better?

Let's make tenants first-class citizens!

# Design goals

**Millions of tenants**
in a single cluster

**Resource isolation**

**Strict separation**
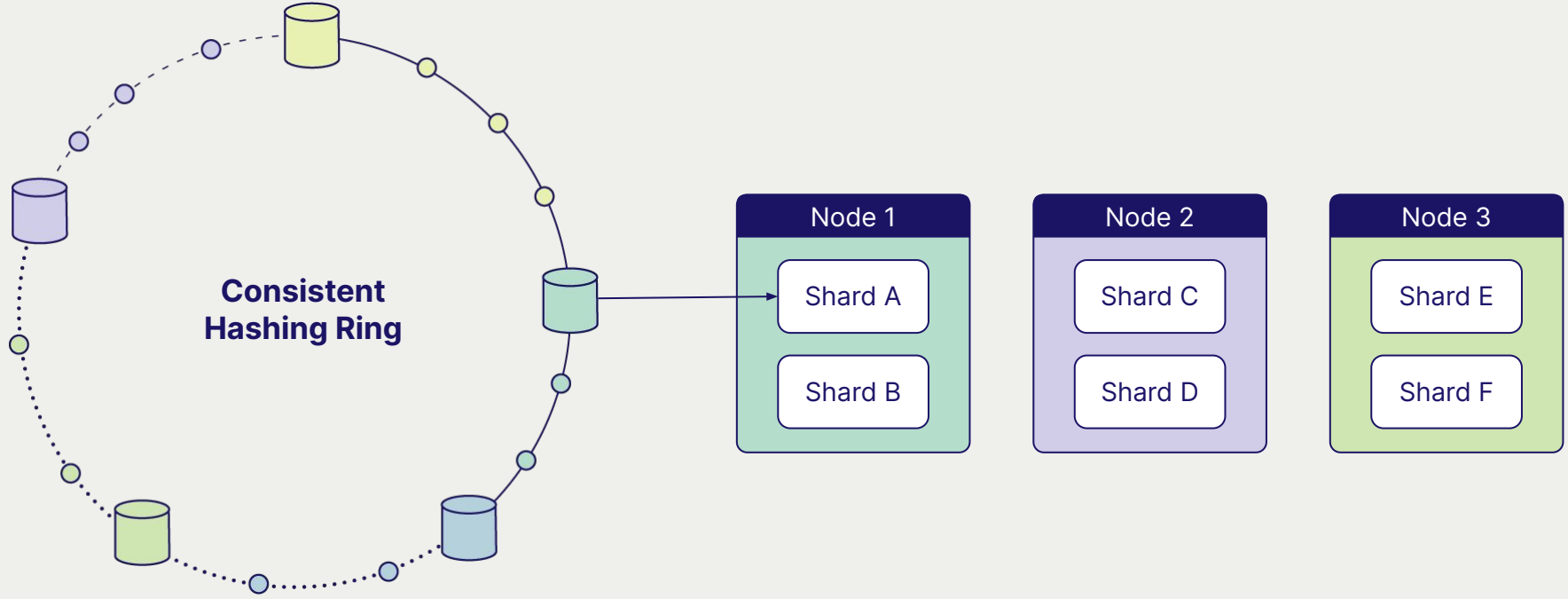for GDPR etc

**Simple offboarding**

**Linear scalability**
Want more tenants?
Add more nodes!

**Only pay for active tenants**

# Sharding would likely play a key role.

# Sharding in a single-tenancy situation



**Consistent Hashing Ring**

**Node 1**
- Shard A
- Shard B

**Node 2**
- Shard C
- Shard D

**Node 3**
- Shard E
- Shard F

What if we took this to the extreme and created a single shard per tenant?

# Pros

**Strict Isolation**
separate storage, separate indexes, resource isolation

**Schema is defined just once**
All shards already belong to the same collection

**Cheap and easy tenant offboardings**
Deleting a tenant is deleting an entire shard (cheap)

# Potential Cons

**Ring–hashing no longer fits**
It's meant to distribute many keys to fewer shards

**Potential for exploding cost**
What is the cost of an empty shard?
What if you have thousands or millions of them?

**Performance Degradation**

SURELY, IT WOULDN'T BE THAT BAD. RIGHT?

RIGHT?

```
runtime: program exceeds 10000-thread limit
fatal error: thread exhaustion

runtime stack:
runtime.throw({0x1053bdbe3?, 0x16b886da0?})
        /usr/local/go/src/runtime/panic.go:1047 +0x40 fp=0x16b886d20 sp=0x16b886cf0 pc=0x104634c00
runtime.checkmcount()
        /usr/local/go/src/runtime/proc.go:789 +0x8c fp=0x16b886d50 sp=0x16b886d20 pc=0x10463861c
runtime.mReserveID()
        /usr/local/go/src/runtime/proc.go:805 +0x3c fp=0x16b886d80 sp=0x16b886d50 pc=0x10463866c
runtime.startm(0x14000064500, 0x0)
        /usr/local/go/src/runtime/proc.go:2403 +0xa8 fp=0x16b886dd0 sp=0x16b886d80 pc=0x10463b6b8
runtime.hand
        /usr
runtime.reta
        /usr
runtime.sysm
        /usr
runtime.msta
        /usr
runtime.msta
        /usr
runtime.msta
        /usr
```

runtime: program exceeds 10000-thread limit
fatal error: thread exhaustion

```
goroutine 1 [semacquire, 1 minutes]:
runtime.gopark(0x106512a20?, 0x0?, 0x0?, 0xe0?, 0x10460c9fc?)
        /usr/local/go/src/runtime/proc.go:381 +0xe0 fp=0x1400053fb00 sp=0x1400053fae0 pc=0x1046375d0
runtime.goparkunlock(...)
        /usr/local/go/src/runtime/proc.go:387
runtime.semacquire1(0x14003da4778, 0xd8?, 0x1, 0x0, 0x98?)
        /usr/local/go/src/runtime/sema.go:160 +0x20c fp=0x1400053fb60 sp=0x1400053fb00 pc=0x104648a9c
sync.runtime_Semacquire(0x140000021a0?)
        /usr/local/go/src/runtime/sema.go:62 +0x2c fp=0x1400053fba0 sp=0x1400053fb60 pc=0x104664b9c
sync.(*WaitGroup).Wait(0x14003da4770)
        /usr/local/go/src/sync/waitgroup.go:116 +0x74 fp=0x1400053fbc0 sp=0x1400053fba0 pc=0x104686b64
github.com/weaviate/weaviate/adapters/handlers/rest.(*Server).Serve(0x14002c21340)
        /Users/etiennedilocker/code/github.com/semi-technologies/weaviate/adapters/handlers/rest/server.go:335 +0x1280 fp=0x140
main.main()
        /Users/etiennedilocker/code/github.com/semi-technologies/weaviate/cmd/weaviate-server/main.go:64 +0x4c8 fp=0x1400053ff7
```
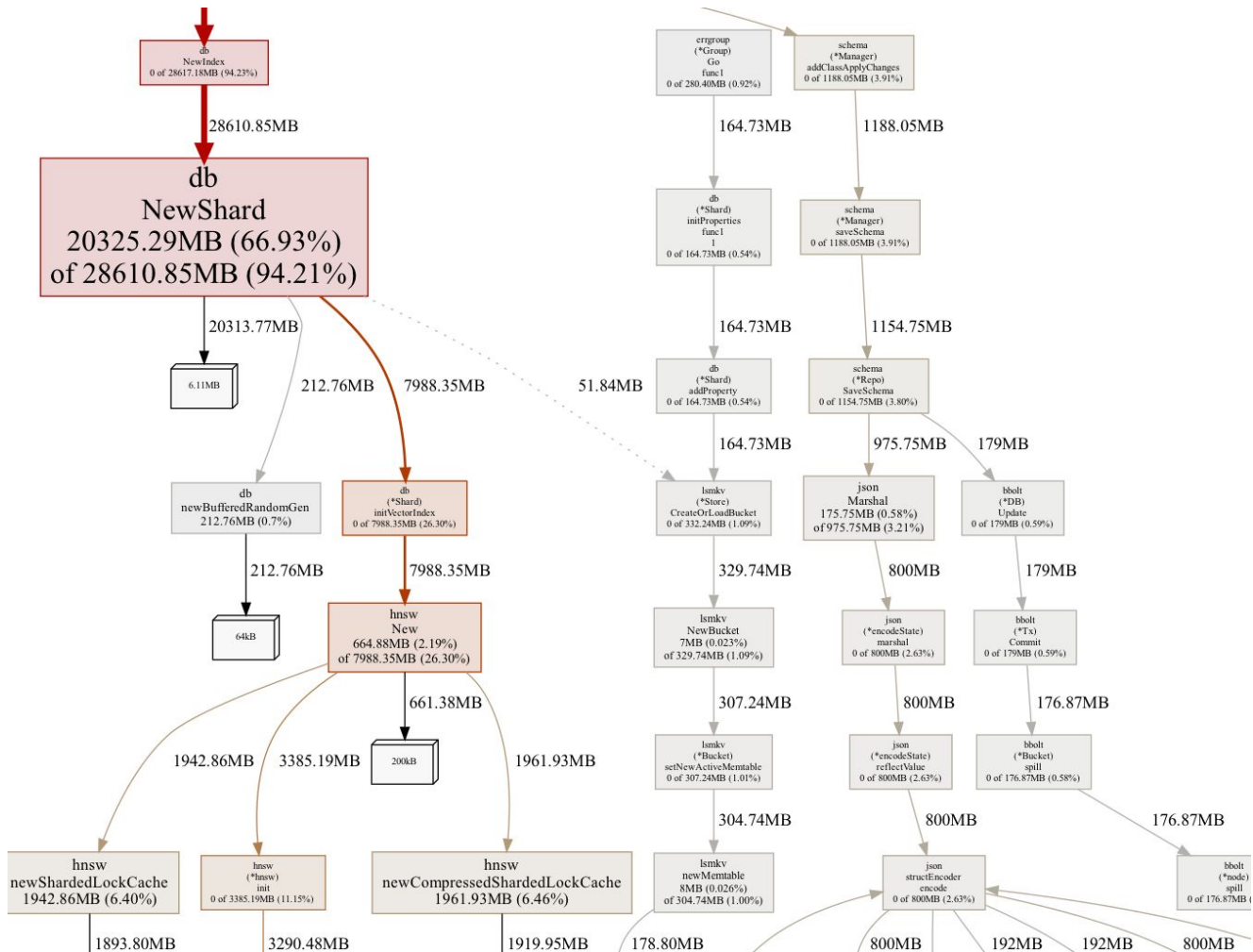
## How expensive would it be?

28GB for 5,000 shards 😱

5.6MB per shard

1M shards → 5.6TB 😅

# Would performance degrade?
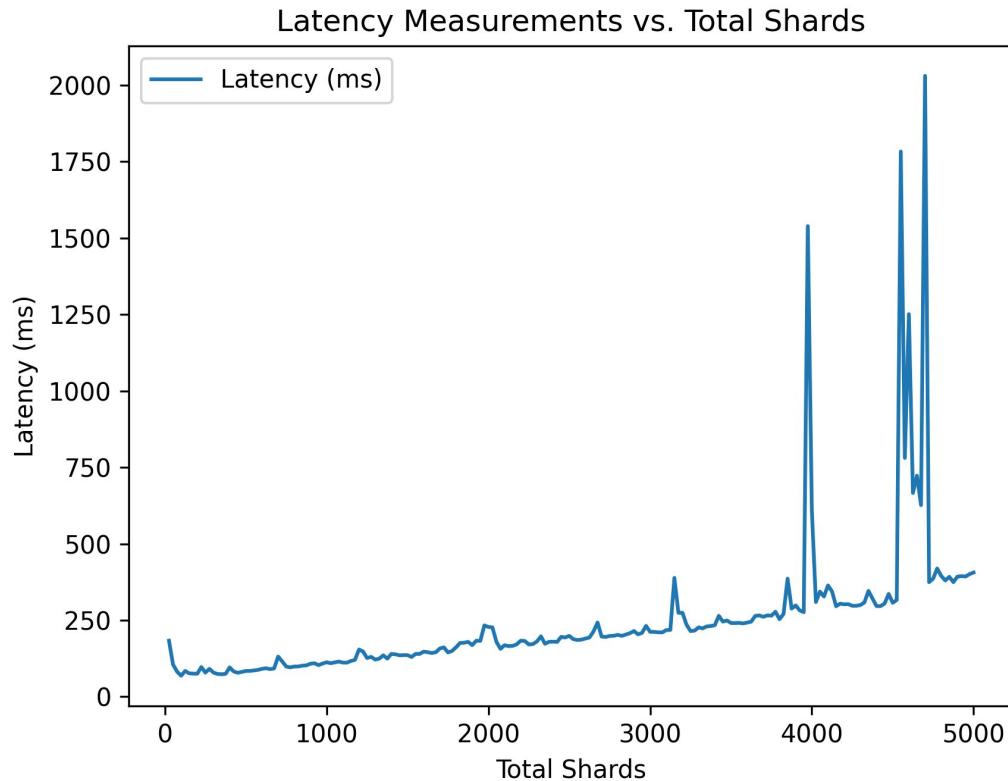
## Creating shards in batches of 25



Latency Measurements vs. Total Shards

# Workarounds don't scale.

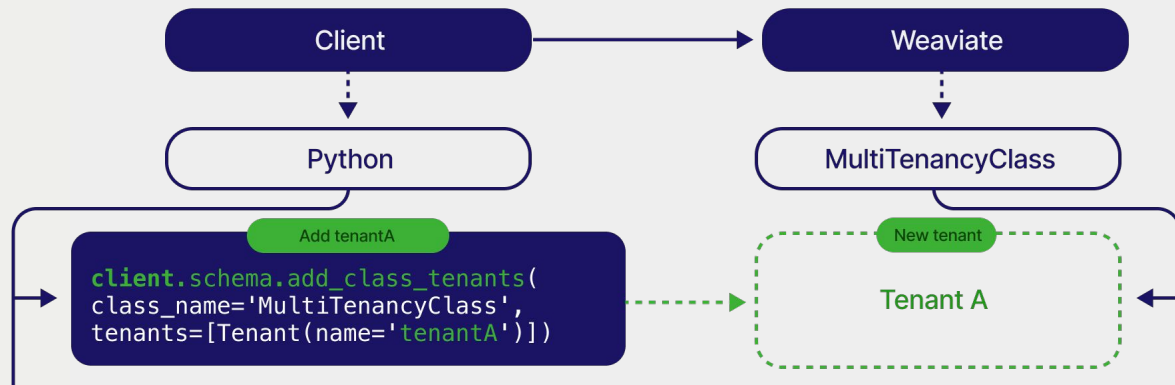These are all addressable problems, but it shows you need an actual Multi-Tenancy solution.

# Simple API

Specify the tenant with each request.

**That's it.**



```
client.schema.add_class_tenants(
class_name='MultiTenancyClass',
tenants=[Tenant(name='tenantA')])
```
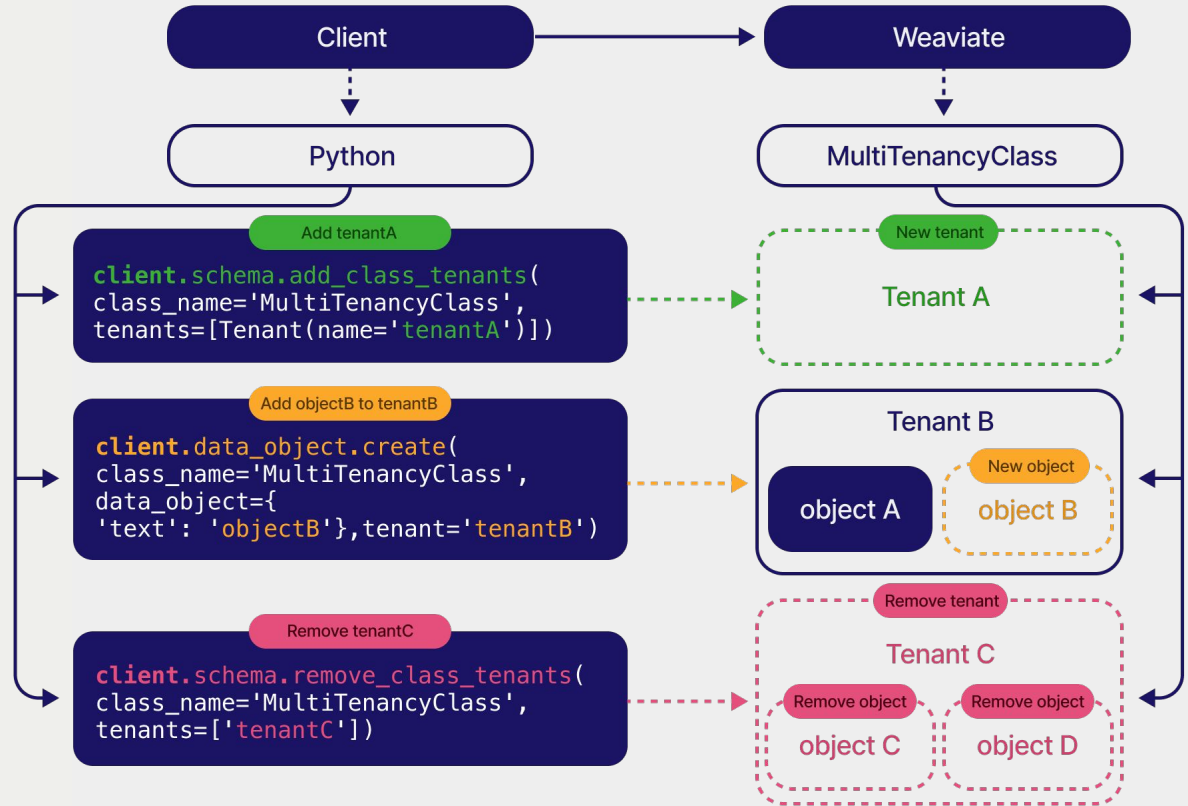
# Simple API

Specify the tenant with each request.

**That's it.**

Client → Weaviate

Python

MultiTenancyClass

**Add tenantA**

```
client.schema.add_class_tenants(
class_name='MultiTenancyClass',
tenants=[Tenant(name='tenantA')])
```

**New tenant**

Tenant A

**Add objectB to tenantB**

```
client.data_object.create(
class_name='MultiTenancyClass',
data_object={
'text': 'objectB'},tenant='tenantB')
```

Tenant B

object A

**New object**

object B

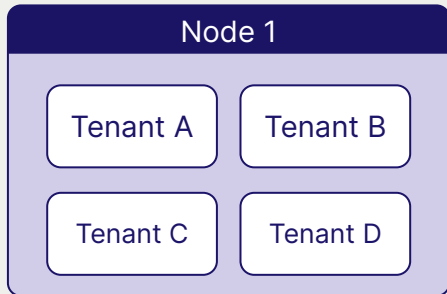# Simple API

Specify the tenant with each request.

**That's it.**

Client → Weaviate

Client ⇢ Python

Weaviate ⇢ MultiTenancyClass

**Add tenantA**
```
client.schema.add_class_tenants(
class_name='MultiTenancyClass',
tenants=[Tenant(name='tenantA')])
```
→ **New tenant** → Tenant A

**Add objectB to tenantB**
```
client.data_object.create(
class_name='MultiTenancyClass',
data_object={
'text': 'objectB'},tenant='tenantB')
```
→ Tenant B — object A, **New object** object B

**Remove tenantC**
```
client.schema.remove_class_tenants(
class_name='MultiTenancyClass',
tenants=['tenantC'])
```
→ **Remove tenant** Tenant C — **Remove object** object C, **Remove object** object D

# Shard routing is a simple (replicated) lookup list

...
Tenant C → Node 1
Tenant D → Node 1
Tenant E → Node 2
Tenant F → Node 2
...

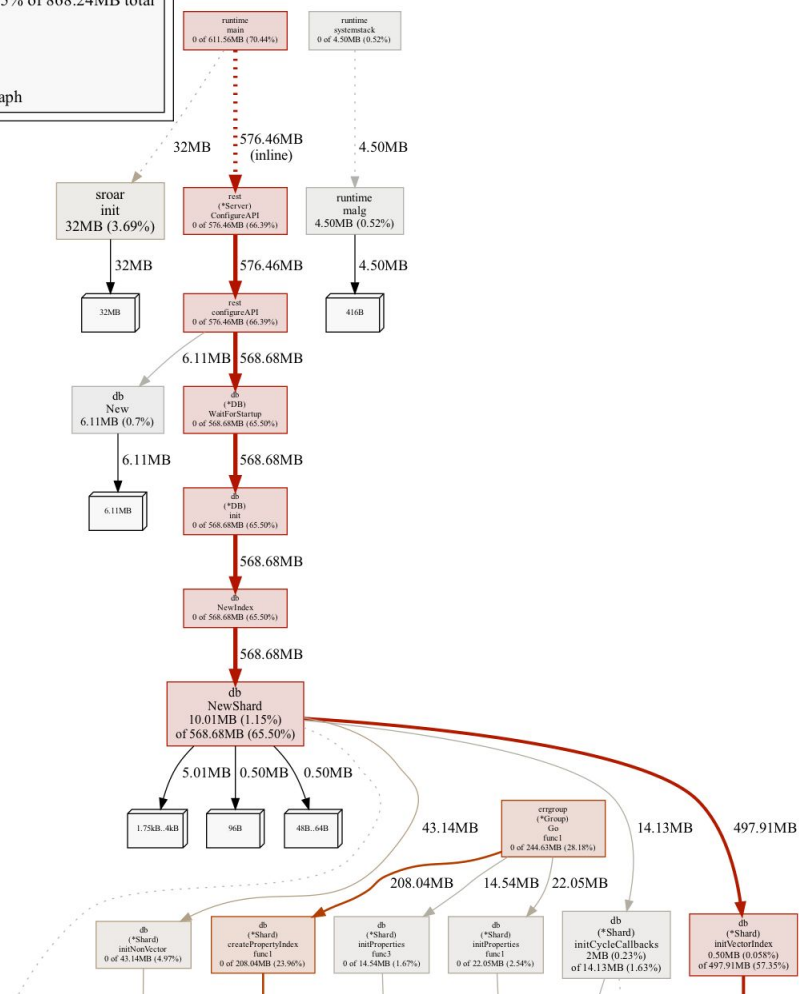| Node 1 | | Node 2 | | Node 3 | |
| --- | --- | --- | --- | --- | --- |
| Tenant A | Tenant B | Tenant E | Tenant F | Tenant I | Tenant J |
| Tenant C | Tenant D | Tenant G | Tenant H | Tenant K | Tenant L |

# 870MB for 5,000 shards

# 174KB per shard

# 1M shards -> 174GB

Type: inuse_space
Time: Sep 20, 2023 at 3:55pm (CEST)
Showing nodes accounting for 827.03MB, 95.25% of 868.24MB total
Dropped 126 nodes (cum <= 4.34MB)
Dropped 4 edges (freq <= 0.87MB)
Showing top 54 nodes out of 62

See https://git.io/JfYMW for how to read the graph

**Performance no longer degrades when adding tenants.**
Creating shards in batches of 25

Latency Measurements vs. Total Shards (Multi-Tenancy v1.20)

Chart uses an identical scale to the previous one (0<x<2000)

**Automatic Multi-Tenancy Load test (Importing)**

## Querying

**QPS (avg over 30s sliding window)**

15000
10000
5000
0

21:20          21:25          21:30

— sum(rate(vector_query_seconds_count[30s]))

**Mean QPS**

15991

**Query Success Rate**

10000

0

21:20          21:25          21:30

— Success
— Failure

**Query success rate**

100.0%

**Currently active users**

3000
2000
1000
0

21:20          21:25          21:30

— sum(querying_users_total)

**Active Users**

2750

**Ready Weaviate Pods**

10

5

0

21:20          21:25          21:30

— sum(kube_pod_status_ready{pod=~"weaviate.*",condition="true"})

**Ready Weaviate pods**

12

**Mean Query latency**

40 ms

20 ms

0 s

21:20          21:25          21:30

— sum(rate(vector_query_seconds_sum[1m0s]))/sum(rate(vector_query_seconds_count[1m0s]))

**Mean Query latency**

35.8 ms
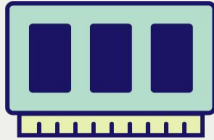
# Automatic Multi-Tenancy Load test (Querying)

# How can Multi-Tenancy reduce cost?

# Storage Tiers and their cost

$$$ fast

**Memory**
very fast, but very expensive

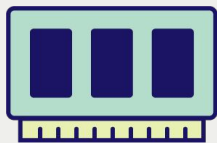**SSD Drives**
medium speed, medium cost

**Cloud Storage**
cheap, but slow

$ slow

# Activating and Deactivating Tenants

**Active Tenants**
are already in memory
and **ready to go**

**Inactive Tenants**
are on disk, but can be loaded
in **a few hundred ms**
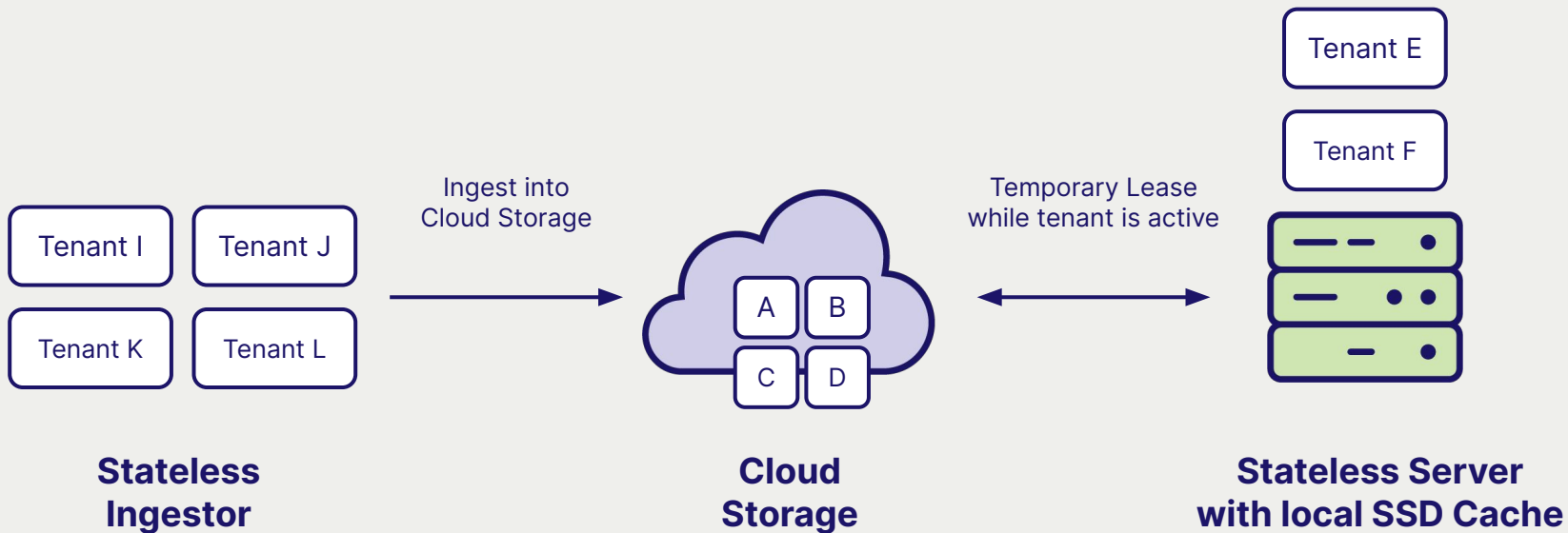
# If 20% of tenants are active at any given time…

If 20% of tenants are active at any given time...

**That's an 80% memory reduction.**

# Separation of Storage & Compute

# Cloud-Storage-based Architecture

**Tenant I**  **Tenant J**

**Tenant K**  **Tenant L**

Ingest into
Cloud Storage

A  B
C  D

Temporary Lease
while tenant is active

**Tenant E**

**Tenant F**

**Stateless
Ingestor**

**Cloud
Storage**

**Stateless Server
with local SSD Cache**

# Recap

Most RAG and search apps require multi-tenancy.

Multi-Tenancy is hard and relying on workaround fails at scale.

With a dedicated Multi-Tenancy solution we can handle the scale and serve tenants efficiently.

There are a lot of new opportunities that couldn't be solved before: Cost reduction, stateless vector dbs, etc.

Thank you

# Connect with us!

🌐 weaviate.io

 weaviate/weaviate

𝕏 @etiennedi

@weaviate_io