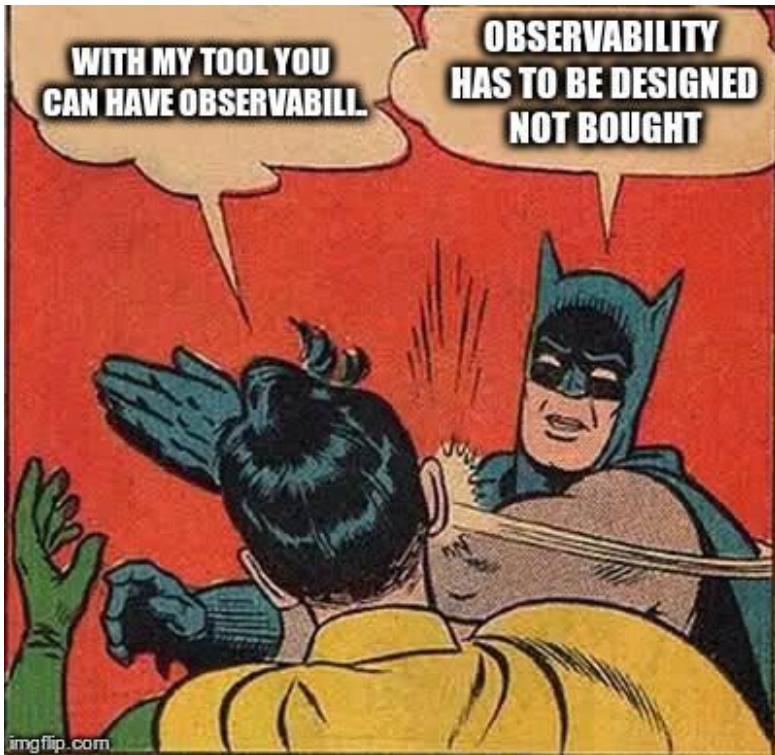


Well Hello There...





imgflip.com

End-to-End Open Source Observability  
Custom-built on **Your Infrastructure**

Build **observability that adapts to your infrastructure**  
Not the other way around.



# The trifecta of 3 OSS tools



Single backend framework  
for capturing telemetry data  
(vendor-agnostic)

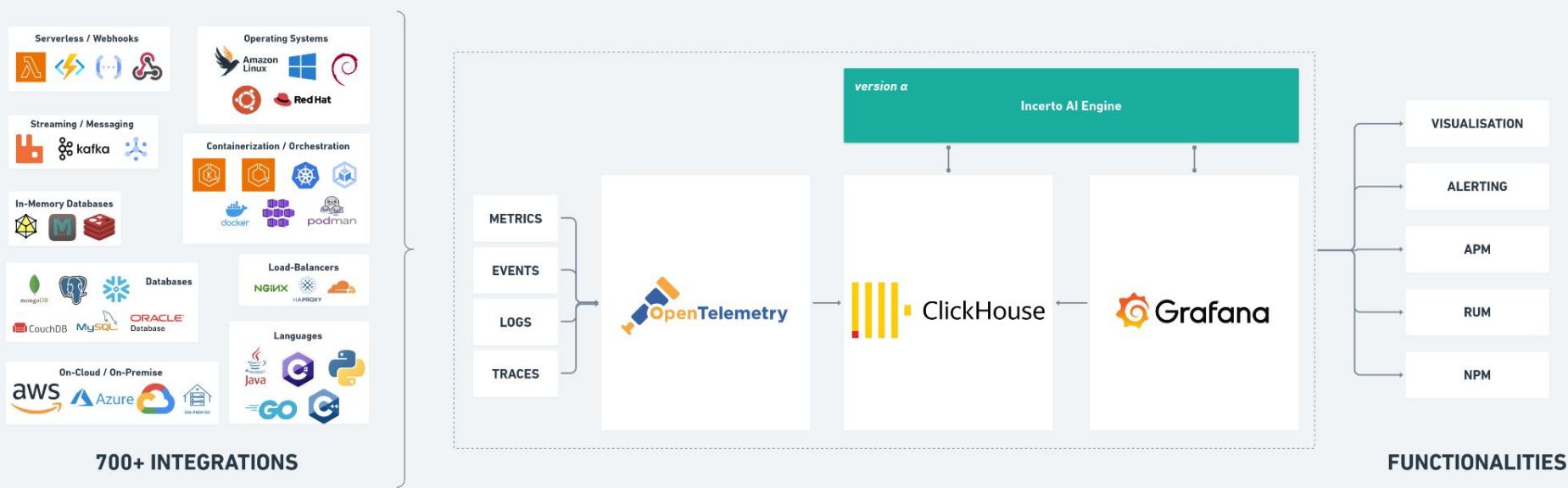


Telemetry data retention



UI for visualization/alerting

# Architecture

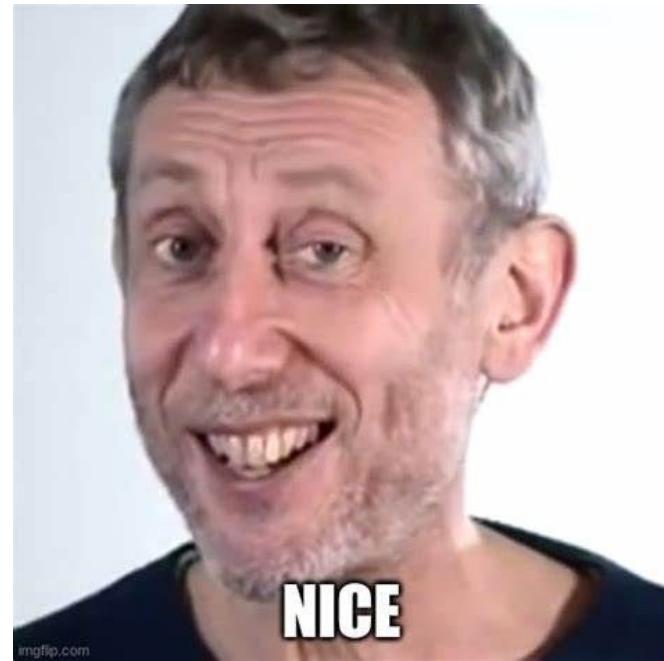


Made with Whimsical

# What have we achieved with Clickhouse ?

**Sub-second Queries** : Includes high cardinality groupings, filtering via various keys, aggregation queries & more.

**Efficient Storage** : Up to 20x compression with S3 storage while juggling TBs of data every month.

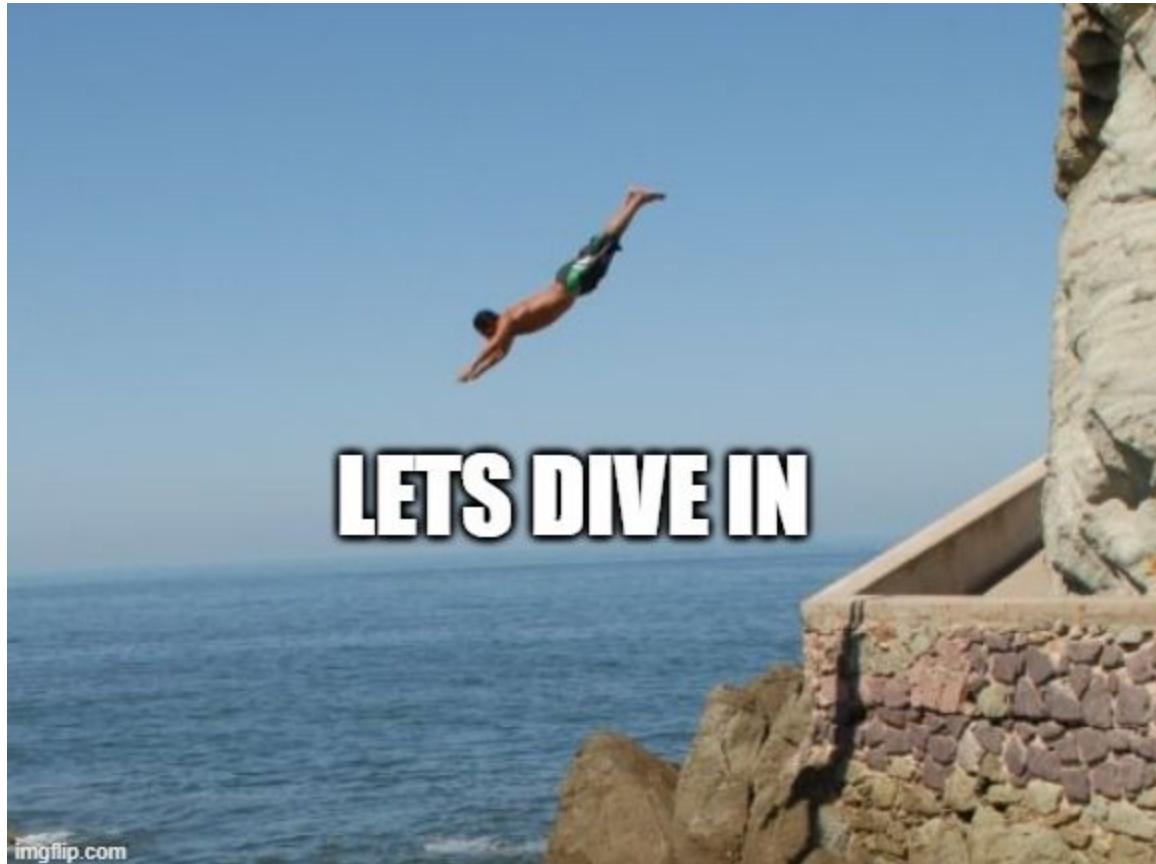


# QUESTION

How to configure Clickhouse for fast queries and efficient storage?

ANSWER

Understanding your query patterns, data and ***somethings about Clickhouse***



**LETS DIVE IN**

imgflip.com

Incerto Technologies

# Explain filtering basics in clickhouse

	URL.bin	UserID.bin	EventTime.bin
row 0	/clck.yandex.ru/file.com ...	1.644.125.792	2014-03-18 09:59:01
row 1	/clck/jsredircnt=1395412 ...	1.550.392.492	2014-03-22 05:55:22
row 2	goal://200906&uinfo=ww-1 ...	2.479.498.648	2014-03-21 05:23:19
	:	:	:
row 8.191	goal://mail.yandex.ru/Ma ...	2.137.667.438	2014-03-21 08:24:50
granule 0 with 8192 rows	→		
row 8.192	goal://mail.yandex.ru/Ma ...	2.137.667.438	2014-03-21 08:25:15
	goal://mail.yandex.ru/Ma ...	2.137.667.438	2014-03-21 08:32:43
	goal://mail.yandex.ru/Ma ...	2.137.667.438	2014-03-21 08:33:02
	:	:	:
row 16.383	goal://r52-echo/#compose ...	1.524.699.296	2014-03-20 05:29:42
granule 1 with 8192 rows	→		
	:	:	:
row 8.863.744	https://wroad.php?show/7 ...	1.878.658.680	2014-03-17 03:55:28
	https://wroad.rt.com.tr/ ...	3.849.470.917	2014-03-20 07:31:39
	https://wroad.rt.com.tr/ ...	3.849.470.917	2014-03-20 07:31:54
	:	:	:
row 8.867.680	res://m.me/politic/stati ...	3.560.941.935	2014-03-19 10:48:46
granule 1082 with 3937 rows	→		

# We want this query to run fast!!

```
SELECT
    timestamp, traceId, spanId, symbolName, logBody,
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE
AND
    traceId = `3b047adb6964601e01688f31b2f0d7a`
AND
    timestamp BETWEEN t1 and t2
```

# Focus on filtering columns

```
SELECT
    timestamp, traceId, spanId, symbol, logBody
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE
AND
    traceId = `3b047adbf6964601e01688f31b2f07a`
AND
    timestamp BETWEEN t1 and t2
```

# Hierarchy and Cardinality

```
SELECT
    timestamp, traceId, spanId, symbol, logBody
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE cardinality = 2
AND
    traceId = `3b047adb6964601e01688f31b2f07a` cardinality > billions
AND
    timestamp BETWEEN t1 and t2 cardinality > billions
```

# Hierarchy and Cardinality

```
SELECT
    timestamp, traceId, spanId, symbol, logBody
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE cardinality = 2
AND
    traceId = `3b047adb6964601e01688f31b2f07a` cardinality > billions
AND
    timestamp BETWEEN t1 and t2 cardinality > billions
```

Rule1 : Increasing order of cardinality

# Hierarchy and Cardinality

```
SELECT
    timestamp, traceId, spanId, symbol, logBody
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE cardinality = 2
AND
    traceId = `3b047adb6964601e01688f31b2f07a` cardinality > billions
AND
    timestamp BETWEEN t1 and t2 cardinality > billions
```

Rule1 : Increasing order of cardinality

Rule2 : Atmost 1 high cardinality key

PrimaryKey := (exchId, traceId)

# Hierarchy and Cardinality

```
SELECT
    timestamp, traceId, spanId, symbol, logBody
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE cardinality = 2
AND
    traceId = `3b047adb6964601e01688f31b2f07a` cardinality > billions
AND
    timestamp BETWEEN t1 and t2 cardinality > billions
```

Rule1 : Increasing order of cardinality

Rule2 : Atmost 1 high cardinality key

PrimaryKey := (exchId, traceId)

# No TraceId

```
SELECT
    timestamp, traceId, spanId, symbol, logBody
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE cardinality = 2
AND
    timestamp BETWEEN t1 and t2 cardinality > billions
```

PrimaryKey := (exchId, traceId)

# Enters Partition Key

PartitionKey := toDate(timestamp)

```
SELECT
    timestamp, traceId, spanId, symbol, logBody
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE cardinality = 2
AND
    timestamp BETWEEN t1 and t2 cardinality > billions
```

# Enters Min Max Skip Indices

```
SELECT
    timestamp, traceId, spanId, symbol, logBody
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE cardinality = 2
AND
    timestamp BETWEEN t1 and t2 cardinality > billions
```

PrimaryKey := (exchId, traceId)

PartitionKey := toDate(timestamp)

addIndex := timestamp(type=Minmax,  
Granularity=3)

# Random High Cardinality

```
SELECT
    timestamp, traceId, spanId, symbol, logBody
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE
    AND timestamp BETWEEN t1 and t2
    AND spanId = fa1152363fac23bf
    cardinality > billions
    cardinality > billions
AND
    cardinality > billions
```

PrimaryKey := (exchId, traceId)

PartitionKey := toDate(timestamp)

addIndex := timestamp(type=Minmax,  
Granularity=3)

# Random High Cardinality

```
SELECT
    timestamp, traceId, spanId, symbol, logBody
FROM
    brokerage.orders_logs
WHERE
    exchID = NSE   cardinality = 2
AND
    timestamp BETWEEN t1 and t2   cardinality > billions
AND
    spanId = fa1152363fac23bf  cardinality > billions
```

PrimaryKey := (exchId, traceId)

PartitionKey := toDate(timestamp)

addIndex := timestamp(type=Minmax,  
Granularity=3)

addIndex := spanId(type=bloom\_filter,  
p=0.01)

A photograph of Michael Scott from the TV show "The Office". He is wearing a dark suit, white shirt, and patterned tie. His hands are clasped under his chin, and he has a slight smile. He is looking directly at the camera with a curious expression. A speech bubble is overlaid on the left side of the image.

GO ON.  
TELL ME MORE

# Quick Recap

**Queries:** List down queries in order of priority. Note down the filtering columns.

**Data:** Know types of columns, number of unique values, range. And for the noted columns only, understand the hierarchy – through grouping, frequencies etc.

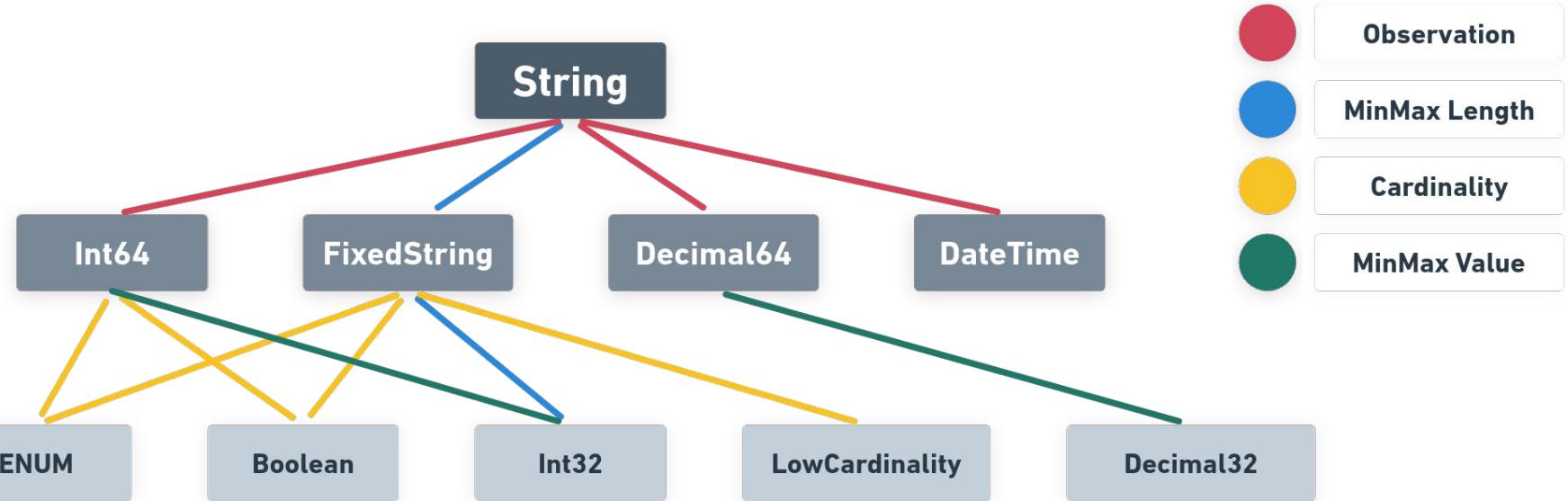
**“Somethings”:** Primary keys, Ordering keys, Skip indexes, Materialized views and Projections.

# Storage Optimisations

How to go from 100GB to 2GB ?



# How to get that 60%?





Use ZSTD(1)

Use Right Primary Key

Use Right Ordering Key

# Top Brokerage firm of India

## Constraints

- **~300GB** logs data/day
- **40TB** telemetry data/month
- Data retention for **7+** years
- **SEBI** compliance specific needs w.r.t auditing
- High **\$\$\$** monthly costs with other vendors



# Result

**“2 SHARD 2 REPLICA”  
CLICKHOUSE CLUSTER**

**4 x (16vCPU 32GB RAM)  
(c7a.4xlarge/c7g.4xlarge)**



# Stage Clickhouse DB snapshot

Database	parts.table	Table								
		rows	latest_modification	disk_size	primary_keys_size	engine	bytes_size	compressed_size	uncompressed_size	
	otel_logs	7423950120	2024-03-13 14:10:...	680.40 GiB	38.39 MiB	ReplicatedMergeT...	730574491252	638.44 GiB	6.55 TiB	
	otel_traces_d_trac...	7545989519	2024-03-13 14:05...	176.98 GiB	39.65 MiB	MergeTree	190027978382	169.69 GiB	343.28 GiB	
	otel_metrics_sum	1471806603	2024-03-13 14:10:...	10.95 GiB	39.30 MiB	ReplicatedMergeT...	11756313207	10.39 GiB	466.80 GiB	
	otel_metrics_gauge	685803259	2024-03-13 14:10:...	6.34 GiB	12.91 MiB	ReplicatedMergeT...	6812333520	5.82 GiB	190.33 GiB	
	otel_traces	23275734	2024-03-13 14:05...	1.43 GiB	272.02 KiB	ReplicatedMergeT...	1532060901	1.16 GiB	8.07 GiB	
	otel_traces_trace_...	36520493	2024-03-13 14:05...	1.02 GiB	321.72 KiB	ReplicatedMergeT...	1100266188	424.52 MiB	717.09 MiB	
	otel_metrics_hist...	42220012	2024-03-13 14:10:...	960.79 MiB	3.58 MiB	ReplicatedMergeT...	1007461740	569.60 MiB	17.82 GiB	
	otel_metrics_sum...	154248	2024-03-13 14:10:...	2.05 MiB	8.22 KiB	ReplicatedMergeT...	2148191	0.00 B	0.00 B	
	otel_metrics_d_sum	223087	2023-12-12 21:33:...	1.51 MiB	4.80 KiB	MergeTree	1586328	1.29 MiB	54.46 MiB	
	otel_metrics_d_ga...	52701	2023-12-12 20:48:...	326.57 KiB	981.00 B	MergeTree	334411	256.08 KiB	10.91 MiB	
	otel_metrics_d_his...	4574	2023-12-12 21:12:...	166.77 KiB	996.00 B	MergeTree	170776	0.00 B	0.00 B	

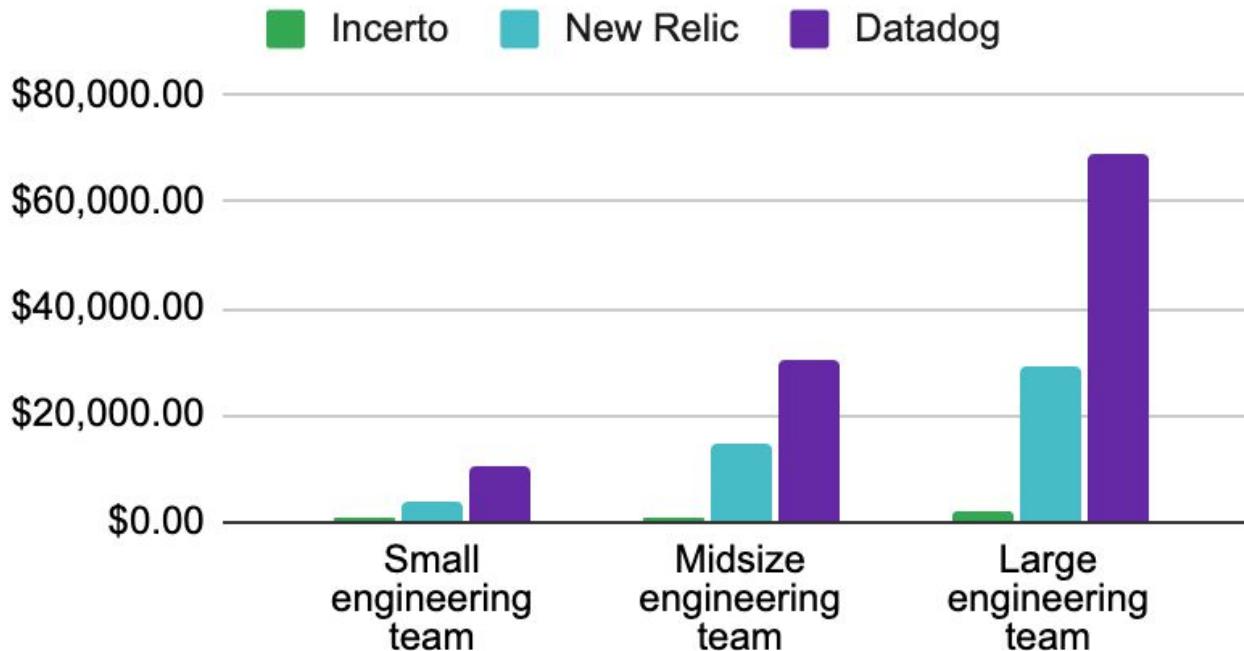
# Retention Cost per GB per month

	A	B	C	D	E
1	Service	1GB	10GB	100GB	1TB
2	Datadog 😢	\$5.080	\$50.805	\$508.047	\$5,080.469
3	Honeycomb	\$2.539	\$25.391	\$253.906	\$2,539.063
4	Logz	\$1.170	\$11.700	\$117.000	\$1,170.000
5	Coralogix	\$1.150	\$11.500	\$115.000	\$1,150.000
6	CW logs	\$0.530	\$5.300	\$53.000	\$530.000
7	Highlight	\$0.50	\$5.00	\$50.00	\$500.00
8	Grafana	\$0.50	\$5.00	\$50.00	\$500.00
9	HyperDX	\$0.40	\$4.00	\$40.00	\$400.00
10	Signoz	\$0.30	\$3.00	\$30.00	\$300.00
11	NewRelic 😢	\$0.300	\$3.000	\$30.000	\$300.000
12	BetterStack	\$0.250	\$2.500	\$25.000	\$250.000
13	Axiom	\$0.150	\$1.500	\$15.000	\$150.000
14	Clickhouse 😎	\$0.047	\$0.471	\$4.710	\$47.100
15	S3 Standard	\$0.023	\$0.230	\$2.300	\$23.000

Ref: <https://bit.kevinslin.com/i/140460591/volume-based-fixed-cost-comparison>

# Incero's stack pricing

## Incero vs New Relic vs Datadog



## Traces &amp; Logs

## Traces

service-A: GET /service-A/items 262.21ms  
2024-03-16 12:01:38.442 GET 200 /check

## Span Filters

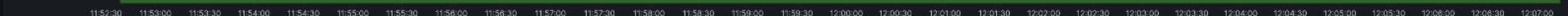
GET /check http send [800.09μs]	800.09μs
GET /check http send [2.44ms]	2.44ms
service-A GET (17.19ms)	17.19ms
service-C GET /service-C/items (13.2ms)	13.2ms
GET (2.99ms)	2.99ms
go-service /region (68.09μs)	68.09μs
getRegion (1.61μs)	1.61μs
service-C GET (2.42ms)	2.42ms
go-service /version (55.3μs)	55.3μs
getVersion (1.25μs)	1.25μs
service-C SELECT (1.71ms)	1.71ms
GET /service-C/items http send (519.8μs)	519.8μs
GET /service-C/items http send (212.6μs)	212.6μs
service-A GET /service-A/items http send (761.31μs)	761.31μs
GET /service-A/items http send (279.3μs)	279.3μs
GET /service-A/items (39.25ms)	39.25ms
GET (6.47ms)	6.47ms
service-B GET /check (2.4ms)	2.4ms
GET /check http send (731.33μs)	731.33μs
GET /check http send (259.97μs)	259.97μs

## Availability

## Error Traces

## Avg Latency

## Uptime Service



## Service Info

ServiceName	errorTraces	errorLogs	ResourceInfo	SpanInfo	LoginInfo	ResourceInfoLog
service-A	265	143	{ "host.name": [ "example-host" ], "os": [ "Windows" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }	{ "http.flavor": [ "1.1" ], "http.host": [ "otelServiceName": "service-A" ], "http.method": "GET", "http.route": "/check", "http.status": 200, "http.time": 1000, "span.id": "1234567890" }	{ "host.name": [ "example-host" ], "os": [ "Windows" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }	{ "host.name": [ "example-host" ], "os": [ "Windows" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }
service-B	202	67	{ "host.name": [ "example-host" ], "os": [ "Linux" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }	{ "http.flavor": [ "1.1" ], "http.host": [ "otelServiceName": "service-B" ], "http.method": "GET", "http.route": "/check", "http.status": 200, "http.time": 1000, "span.id": "1234567890" }	{ "host.name": [ "example-host" ], "os": [ "Linux" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }	{ "host.name": [ "example-host" ], "os": [ "Linux" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }
service-C	119	45	{ "host.name": [ "example-host" ], "os": [ "Linux" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }	{ "http.flavor": [ "1.1" ], "http.host": [ "otelServiceName": "service-C" ], "http.method": "POST", "http.route": "/items", "http.status": 200, "http.time": 1000, "span.id": "1234567890" }	{ "host.name": [ "example-host" ], "os": [ "Linux" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }	{ "host.name": [ "example-host" ], "os": [ "Linux" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }
go-service	43	0	{ "host.name": [ "example-host" ], "os": [ "Linux" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }	{ "http.method": [ "GET" ], "http.route": [ "/check" ], "http.status": 200, "span.id": "1234567890" }	{ "host.name": [ "example-host" ], "os": [ "Linux" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }	{ "host.name": [ "example-host" ], "os": [ "Linux" ], "cpu": 10, "mem": 1000, "disk": 1000, "network": { "in": 100, "out": 100 } }

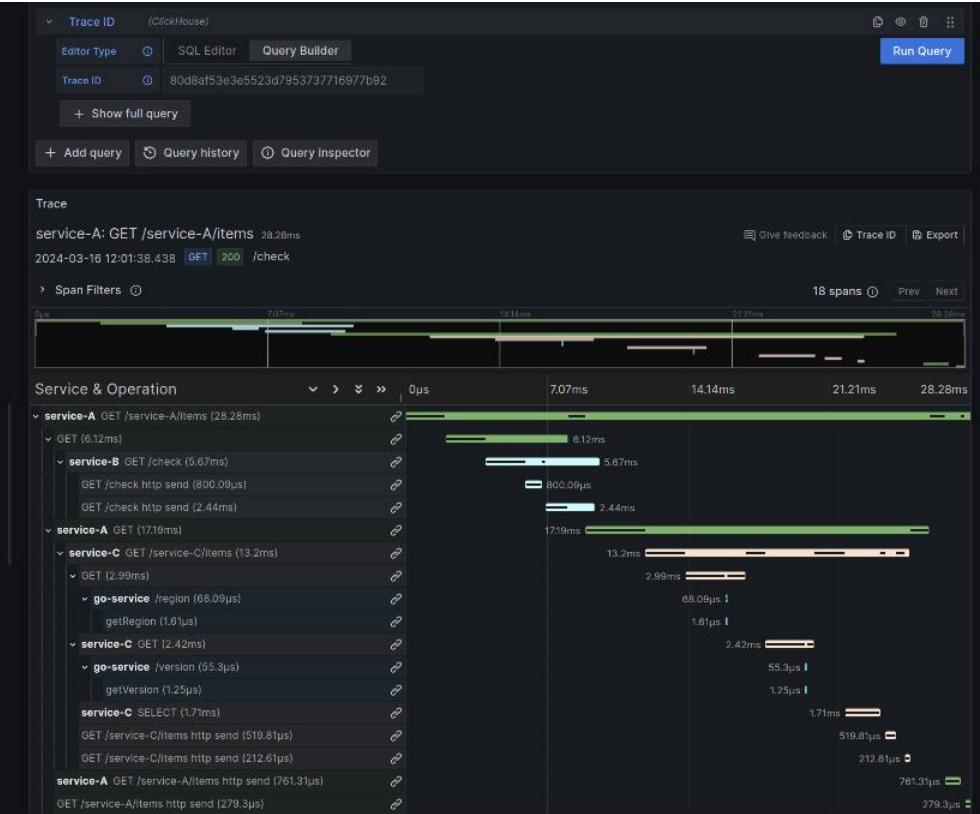
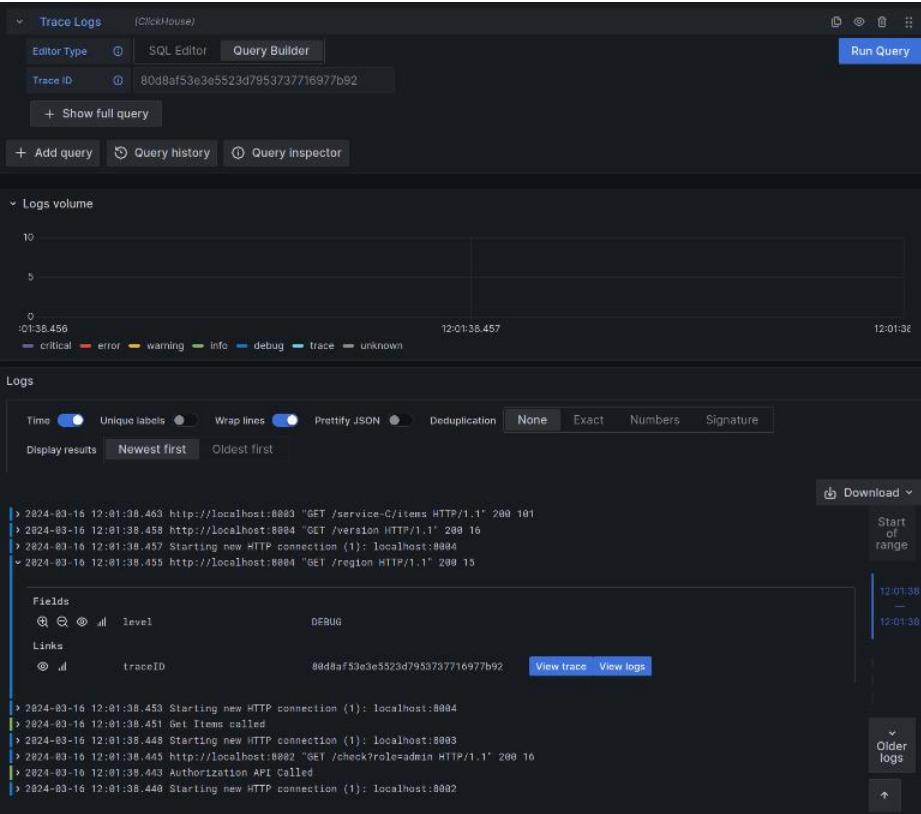
## All Logs

```
> 2024-03-16 12:01:38.702 http://localhost:8003 'POST /service-C/items HTTP/1.1' 200 50
> 2024-03-16 12:01:38.944 Add item called
> 2024-03-16 12:01:38.991 Starting new HTTP connection () localhost:8003
> 2024-03-16 12:01:38.989 http://localhost:8002 'GET /check?role=admin HTTP/1.1' 200 16
> 2024-03-16 12:01:39.085 Authorization API Called
> 2024-03-16 12:01:39.085 Starting new HTTP connection () localhost:8002
> 2024-03-16 12:01:39.108 Error: Price can't be less than 100
> 2024-03-16 12:01:39.465 Error: Price can't be less than 100
> 2024-03-16 12:01:39.453 http://localhost:8003 'POST /service-C/items HTTP/1.1' 200 50
> 2024-03-16 12:01:39.401 Add item called
> 2024-03-16 12:01:39.401 Starting new HTTP connection () localhost:8003
> 2024-03-16 12:01:39.403 http://localhost:8002 'GET /check?role=admin HTTP/1.1' 200 16
> 2024-03-16 12:01:39.485 Authorization API Called
> 2024-03-16 12:01:39.482 Starting new HTTP connection () localhost:8002
> 2024-03-16 12:01:39.615 Error: Price can't be less than 100
```

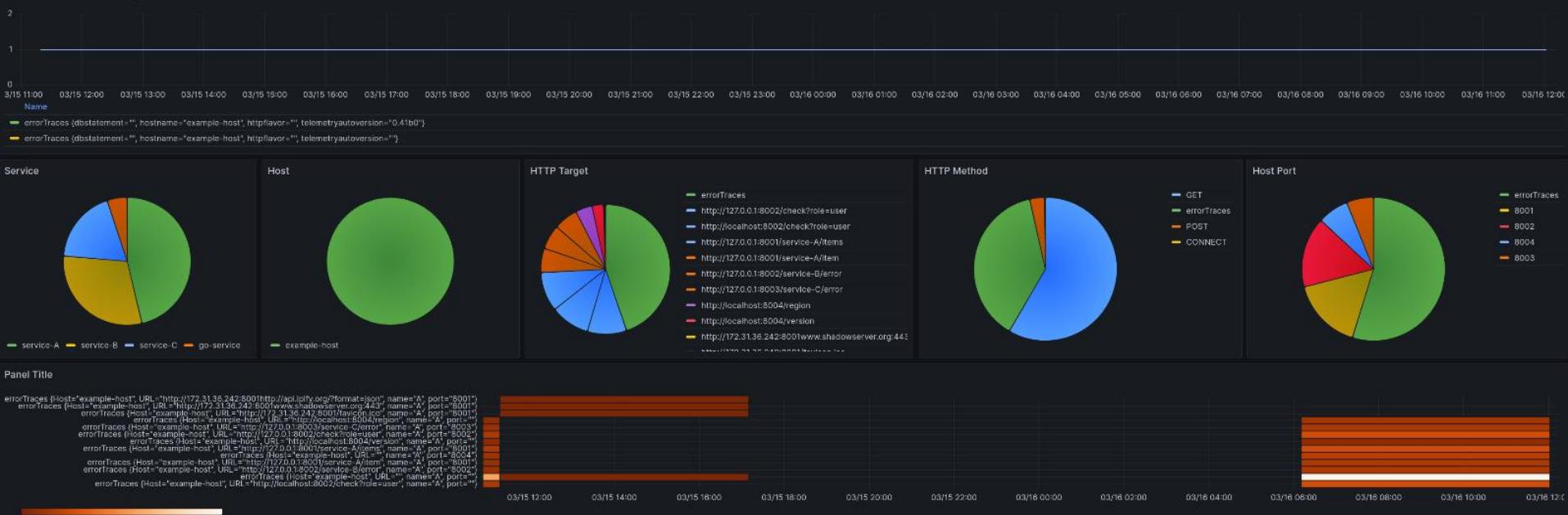
## Service Graph



Legend: mainstat (green), secondarystat (blue), arc\_error (red), arc\_ok (green)



## Dynamic TimeSeries Grouping



Panel Title

Panel Title

service-A: GET / 1d ln

2024-03-15 11:17:23.299 GET 200 /

Give feedback  Trace ID  Export

> Span Filters

158 spans  Prev Next



Grouped Latencies



Slow APIs



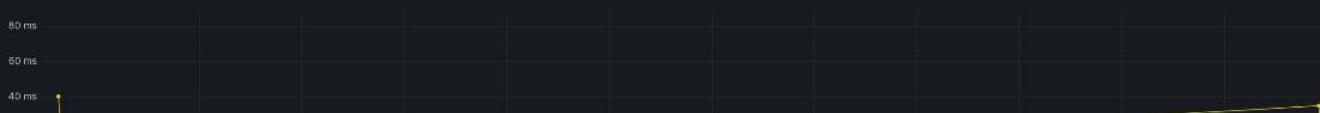
Grouped Latencies



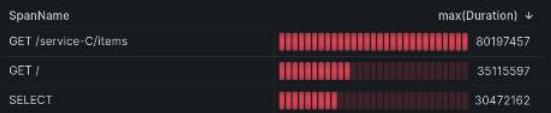
Slow APIs



Grouped Latencies

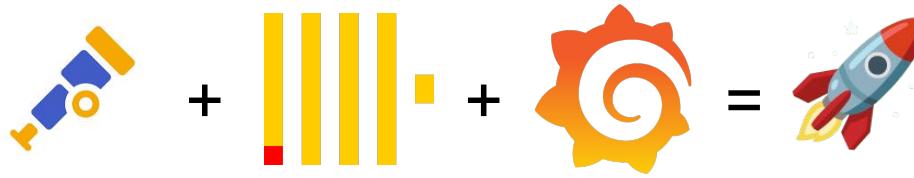


Slow APIs



The key thing  
we want you take home today is

You can ABSOLUTELY build your own  
observability



Or you can just come to us 😎



**Shiva Pundir**

**Co-Founder  
Inceto Technologies**



**Shiva Pundir**

Co-Founder @ Incerto | Building  
Custom Observability for Enterprise

