# Vector Search in ClickHouse

Do you really need a Vector database?

Dale McDiarmid - Product@ClickHouse

**May, 2023**

ClickHouse

# A *very* brief history of <u>search</u>
(more [here](#))

ClickHouse

# Bag of Words + Inverted Indices



Documents

**doc_1**

The moonlight danced on the surface of the calm lake as the night creatures began to stir.

**doc_2**

The sound of rain splashing against the lake was a comforting lullaby for the sleepless night.

**doc_3**

As the rain pattered softly on the surface of the tranquil lake under the glow of the moonlight, the night came alive with the soothing sounds of nature.

Tokenization

(split on whitespace ->
remove punctuation ->
lowercase ->
remove plurals)

## Bag of Words

the
moonlight
danced
on
the
surface
of
the
calm
lake
as
the
night
creatures
began
to
stir
the
sound
of
rain
splashing
against
the
lake
was
a
comforting
lullaby
for
the
sleepless
night

as
the
rain
pattered
softly
on
the
surface
of
the
tranquil
lake
under
the
glow
of
the
moonlight
the
night
came
alive
with
the
soothing
sounds
of
nature

terms to index

## Term Index

| term | frequency | doc ids (doc_id: frequency) |
|------|-----------|------------------------------|
| the | 14 | [ (doc_1, 4), (doc_2, 3), (doc_3, 7) } |
| moonlight | 2 | [ (doc_1, 1), (doc_3, 1) } |
| danced | 1 | [ (doc_1, 1) ] |
| on | 2 | [ (doc_1, 1), (doc_3, 1) ] |
| surface | 2 | [ (doc_1, 1), (doc_3, 1) ] |
| of | 5 | [ (doc_1, 1), (doc_2, 1), (doc_3, 3) ] |
| calm | 1 | [ (doc_1, 1) ] |
| lake | 3 | [ (doc_1, 1), (doc_2, 1), (doc_3, 1) ] |
| as | 2 | [ (doc_1, 1), (doc_2, 1) ] |
| night | 3 | [ (doc_1, 1), (doc_2, 1), (doc_3, 1) ] |
| sound | 2 | [ (doc_2, 1), (doc_3, 1) ] |
| rain | 2 | [ (doc_2, 1), (doc_3, 1) ] |

etc...

# Scoring - Measuring Relevancy - TF-IDF/BM25



TF is the frequency
of a term in a document

$$TF - IDF = TF(t, d) \times IDF(t)$$

IDF - The inverse document frequency
the ratio of documents that include the term

## Term Index

| term | frequency | doc ids (doc_id: frequency) |
|------|-----------|------------------------------|
| the | 14 | [ (doc_1, 4), (doc_2, 3), (doc_3, 7) } |
| moonlight | 2 | [ (doc_1, 1), (doc_3, 1) } |
| danced | 1 | [ (doc_1, 1) ] |
| on | 2 | [ (doc_1, 1), (doc_3, 1) ] |
| surface | 2 | [ (doc_1, 1), (doc_3, 1) ] |
| of | 5 | [ (doc_1, 1), (doc_2, 1), (doc_3, 3) ] |
| calm | 1 | [ (doc_1, 1) ] |
| lake | 3 | [ (doc_1, 1), (doc_2, 1), (doc_3, 1) ] |
| as | 2 | [ (doc_1, 1), (doc_2, 1) ] |
| night | 3 | [ (doc_1, 1), (doc_2, 1), (doc_3, 1) ] |
| sound | 2 | [ (doc_2, 1), (doc_3, 1) ] |
| rain | 2 | [ (doc_2, 1), (doc_3, 1) ] |

etc...

Search: "The sound of the lake"    (assume a logical AND)

### doc_2

| term | TF | IDF | TF x IDF |
|------|-----|------|----------|
| the | 3 | 1/14 | 3/14 |
| sound | 1 | 1/2 | 1/2 |
| of | 1 | 1/5 | 1/5 |
| the | 3 | 1/14 | 3/14 |
| lake | 1 | 1/3 | 1/3 |
| Total (score) | | | 1.46 |

### doc_3

| term | TF | IDF | TF x IDF |
|------|-----|------|----------|
| the | 7 | 1/14 | 7/14 |
| sound | 1 | 1/2 | 1/2 |
| of | 3 | 1/5 | 3/5 |
| the | 7 | 1/14 | 7/14 |
| lake | 1 | 1/3 | 1/3 |
| Total (score) | | | 2.43 |

doc_3 is more relevant

# What is a vector/embedding?

**(more than just text)**

ClickHouse

# Vector vs Embedding

Example 10-dimensional vector

$\langle$ 0.0318, 0.0356, 0.0693 -0.0147, -0.0417, 0.0014, -0.005, -0.0034, -0.0683, -0.0432 $\rangle$
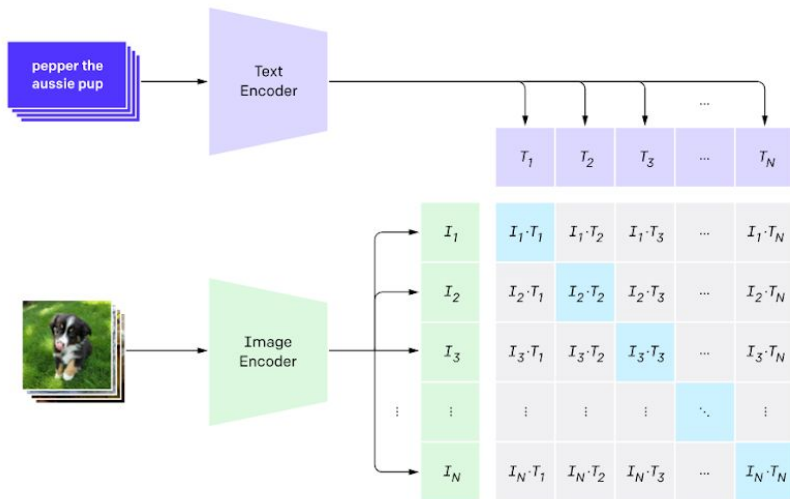
Embedding for the word "moonlight"
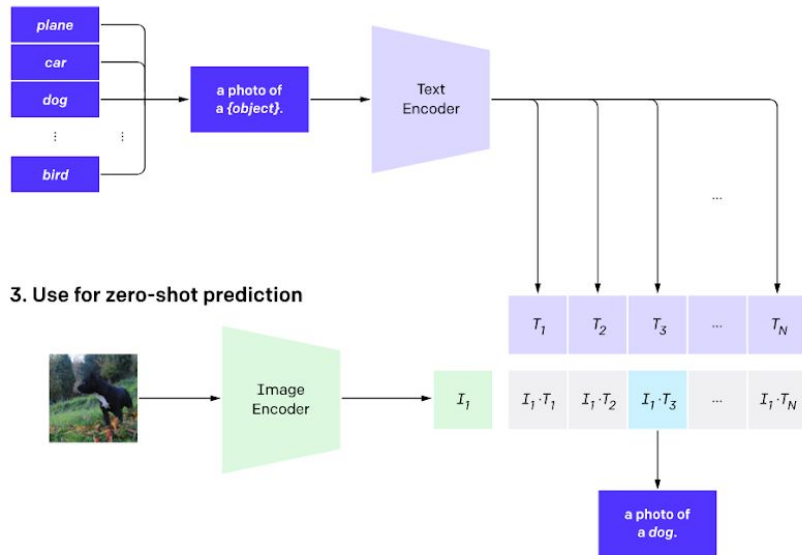
$\langle$ 0.0318, 0.0356, 0.0693 0.567, -0.0417, 0.9675 $\rangle$

$\downarrow$

$\langle$ color temperature, directionality, mood, brightness, warmth, nighttime $\rangle$
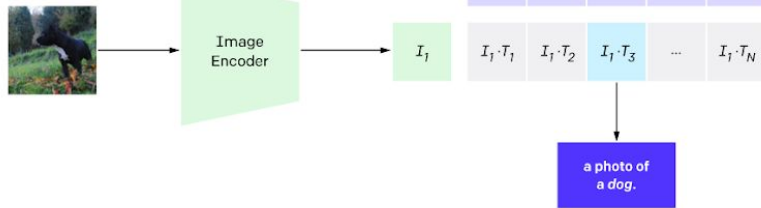
# Producing Embeddings - a lot more than text
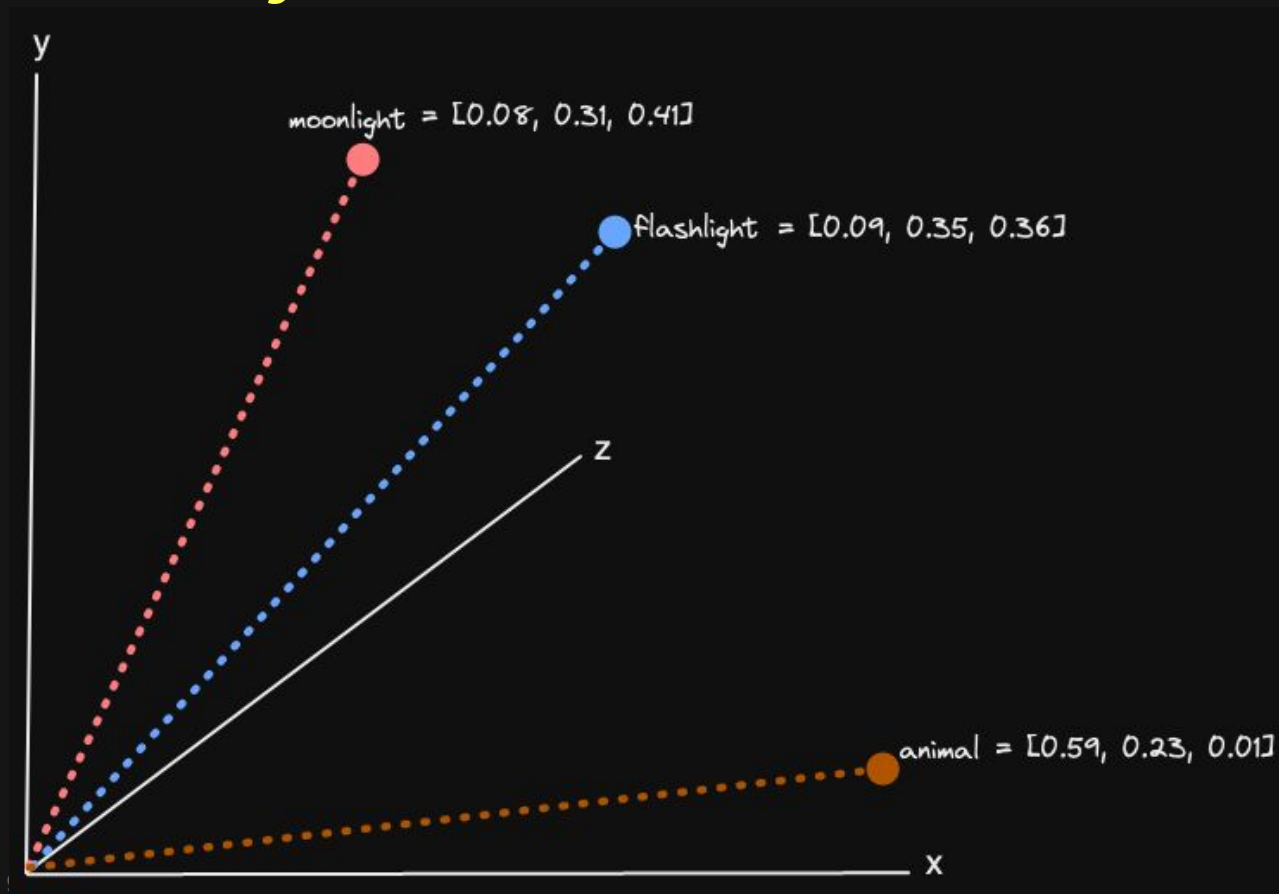


Credit: CLIP -
https://openai.com/research/clip

# The importance of distance
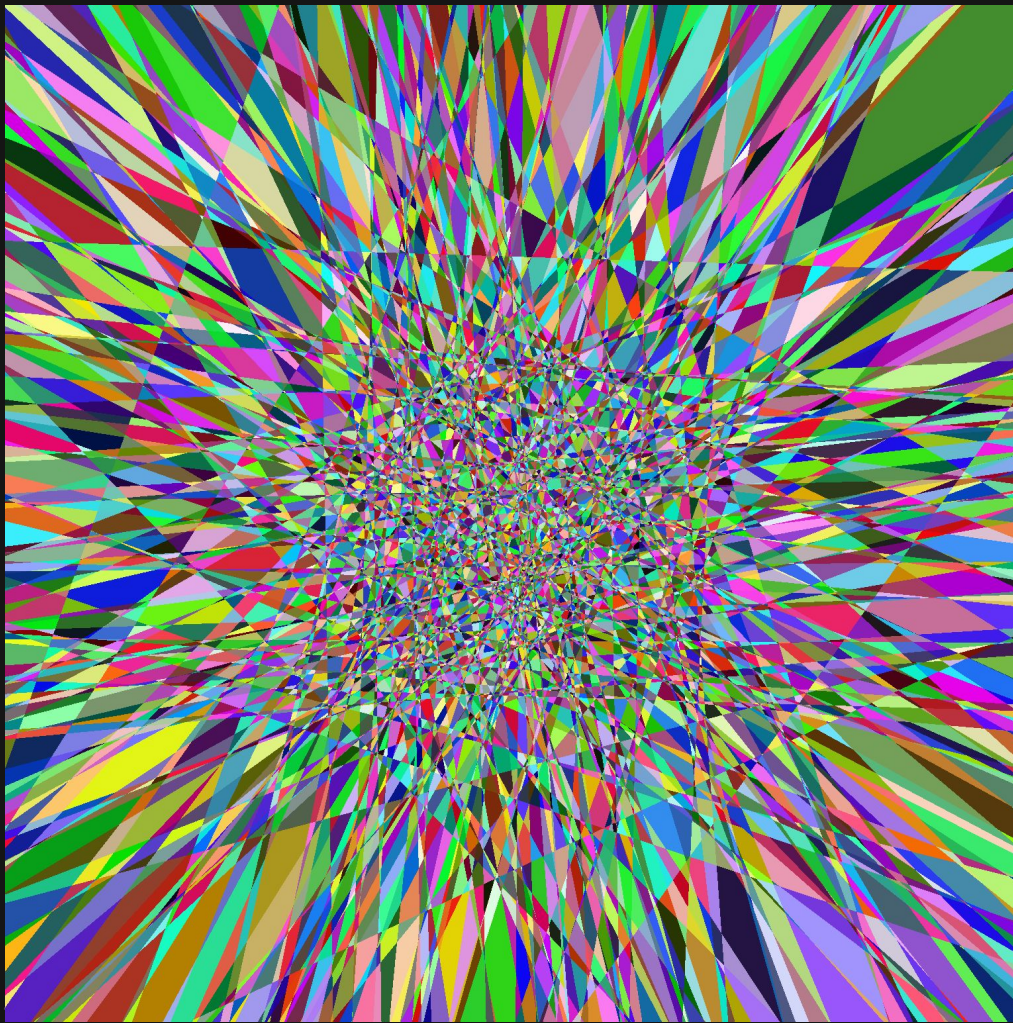
# Similarity measured with distance

# Searching vectors

# Exact linear scan

# Approximate Techniques

- Approximate Nearest Neighbour (ANN)
- Trade accuracy for performance
- Image is for Annoy algorithm
- Annoy builds a tree structure by partitioning the space into smaller and smaller hyperplanes for recursive evaluation

# What problems can we solve?

ClickHouse

# Possible Applications

**01**
### Recommendations
Relevant to e-commerce websites. Beyond embedding text meaning into vectors, page views and past purchases can be encoded.

**02**
### Question Answering
Equivalent meanings can, however, be encoded with vectors that are close, e.g., X and Y.

**03**
### Image and video search
Search for images and videos based on text and vise versa

**04**
### Fraud Detection
Find similar or dissimilar transactions by encoding users' behaviors or log-in patterns into vectors. These can be anomalous behaviors and prevent fraud.

**05**
### Providing Context
Leveraging vector search to provide contextual content to chat applications powered by APIs like ChatGPT.

**06**
### Multilingual search
Allow cross-language searching with the same concept in two languages encoded to the same vector.

# Vector Search in ClickHouse

ClickHouse

- **Vectors are just Array(Float32) - 1.5x to 2x compression**

- **Matching is just a distance function - mainly Euclidean or Cosine but others supported**

- **Full SQL support and aggregations**

- **Approximate Nearest Neighbour supported through Annoy. Used to optimize distance functions.\***

- **UDFs are useful for integrating embedding generation**

**\*Experimental.**

ClickHouse

```
SELECT

    url,

    caption,

    L2Distance(image_embedding, [0.5736801028251648...0.2516217529773712]) AS score

FROM laion WHERE similarity >= 0.2

ORDER BY score ASC

LIMIT 2

FORMAT Vertical


10 rows in set. Elapsed: 1.603 sec. Processed 10.00 million rows, 32.70 GB (6.24 million rows/s.,
20.40 GB/s.)
```

# When to use ClickHouse for Vector Search?

ClickHouse

# ClickHouse for Vector Search

You have/need:

- High performance and scalable linear matching using many cores
- Metadata with your vectors for which you'll benefit from high compression and extremely fast querying/filtering
- Full SQL support in your querying inc. Joins and aggregations
- A very large dataset and do not wish to be memory bound or a small dataset where linear matching is fine (~millions)
- Data in ClickHouse and wish to complement with vectors
- An existing pipeline for generating vectors from a model and don't need tight coupling

# Demo

ClickHouse

# LAION dataset

- [5 billion embeddings](#) created with performance analysis in mind
- Created using a [CLIP (Contrastive Language–Image Pre-training)](#) multi-modal model. Allows images to be classified with a relevant category for previously unknown images not seen during the training process
- Trained on 400m image/caption pairs
- Each row has image metadata and a text + image embedding (768 dimensions)
- Allows searching for images with text and vise versa using Euclidean distance for matching
- We used a [2.2b english subset](#) (4 billion vectors in total)
- We combined embeddings with metadata and made available in Parquet*

# bfloat16

## Floating Point Formats

**bfloat16: Brain Floating Point Format**

Range: ~$1e^{-38}$ to ~$3e^{38}$

Exponent: 8 bits | Mantissa (Significand): 7 bits

| S | E E E E E E E E | M M M M M M M |

**fp32: Single-precision IEEE Floating Point Format**

Range: ~$1e^{-38}$ to ~$3e^{38}$

Exponent: 8 bits | Mantissa (Significand): 23 bits

| S | E E E E E E E E | M M M M M M M M M M M M M M M M M M M M M M M |

**fp16: Half-precision IEEE Floating Point Format**

Range: ~$5.96e^{-8}$ to 65504

Exponent: 5 bits | Mantissa (Significand): 10 bits

| S | E E E E E | M M M M M M M M M M |

```sql
INSERT INTO laion_bfloat16
SELECT
    arrayMap(x -> reinterpretAsFloat32(bitAnd(reinterpretAsUInt32(x),
4294901760)), image_embedding) AS image_embedding,
    arrayMap(x -> reinterpretAsFloat32(bitAnd(reinterpretAsUInt32(x),
4294901760)), text_embedding) AS text_embedding
FROM laion
```

**END**

# Speakers

**Aaron Katz**

CEO @ ClickHouse

**Alexey Milovidov**

CTO @ClickHouse

**Yury  Izrailevsky**

President and
VP of Engineering

# 01 Section title with number

**Author or subtitle**

◧ ClickHouse

# Content heavy slide (2 columns)

# Content slide - less text heavy

# What is ClickHouse?

## Open source

- Developed since 2009
- OSS 2016
- 28k+ Github stars
- 1k+ contributors
- 300+ releases

## column-oriented

- Best for aggregations
- Files per column
- Sorting and indexing
- Background merges

## distributed

- Replication
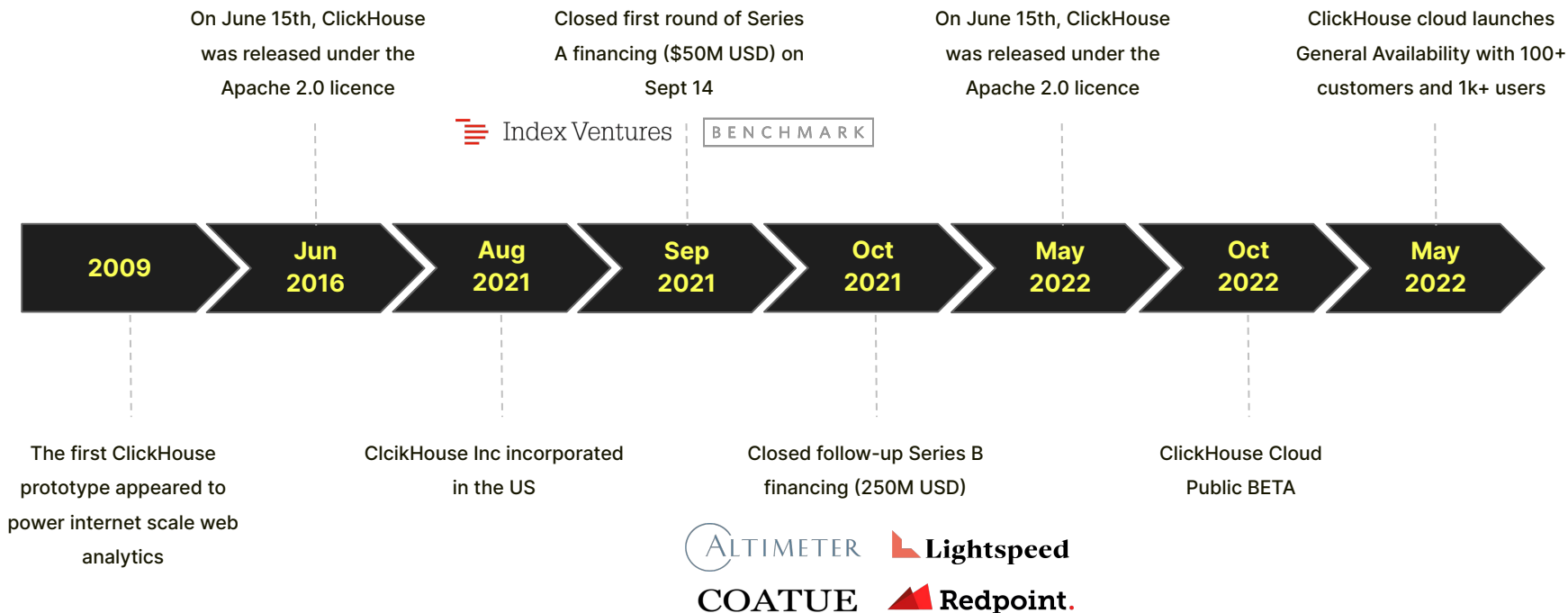- Sharding
- Multi-master
- Cross-region

## OLAP database

- Analytics use cases
- Aggregations
- Visualizations
- Mostly immutable data

# What is ClickHouse? (highlighted)

**Open source**

**column-oriented**

**distributed**

**OLAP database**

| Open source | column-oriented | distributed | OLAP database |
|---|---|---|---|
| Developed since 2009 | Best for aggregations | Replication | Analytics use cases |
| OSS 2016 | Files per column | Sharding | Aggregations |
| 28k+ Github stars | Sorting and indexing | Multi-master | Visualizations |
| 1k+ contributors | Background merges | Cross-region | Mostly immutable data |
| 300+ releases | | | |

# ClickHouse Journey

On June 15th, ClickHouse was released under the Apache 2.0 licence

Closed first round of Series A financing ($50M USD) on Sept 14

Index Ventures   BENCHMARK

On June 15th, ClickHouse was released under the Apache 2.0 licence

ClickHouse cloud launches General Availability with 100+ customers and 1k+ users

| 2009 | Jun 2016 | Aug 2021 | Sep 2021 | Oct 2021 | May 2022 | Oct 2022 | May 2022 |

The first ClickHouse prototype appeared to power internet scale web analytics

ClcikHouse Inc incorporated in the US

Closed follow-up Series B financing (250M USD)

ALTIMETER   Lightspeed
COATUE   Redpoint.

ClickHouse Cloud Public BETA

31

# Feature highlight

- list item 1
- list item 2
- list item 3

# Two cards highlight

## Title here

**Subtitle or a brief description**

- Another bullet here
- Another bullet here
- Another bullet here
- Another bullet here

## Title here

**Subtitle or a brief description**

- Another bullet here
- Another bullet here
- Another bullet here
- Another bullet here

# Three cards highlight

## Title here

**Subtitle or a brief description**

- Another bullet here
- Another bullet here
- Another bullet here
- Another bullet here

## Title here

**Subtitle or a brief description**

- Another bullet here
- Another bullet here
- Another bullet here
- Another bullet here

## Title here

**Subtitle or a brief description**

- Another bullet here
- Another bullet here
- Another bullet here
- Another bullet here

# Use this slide for code (dark)

```
SELECT

    toStartOfMonth(upload_date) AS month,

    sum(view_count) AS `Youtube Views`,

    bar(sum(has_subtitles) / count(), 0.55, 0.7, 100) AS `% Subtitles`

FROM youtube

WHERE (month >= '2020-08-01') AND (month <= '2021-08-01')

GROUP BY month

ORDER BY month ASC


13 rows in set. Elapsed: 0.823 sec Processed 1.07 billion rows, 11.75 GB (1.30 billion rows/s., 14.27 GB/s.)
```

# ClickHouse Cloud - from 0 to 1 in under a year

## Private preview
### May -July 2022

## Public Beta
### Oct 2022

## General Availability
### Dec 2022

**Serverless hosted ClickHouse for key design partners**

- AWS with 3 regions initially

- Limited scalability within predefined limits

- Basic cloud console for ClickHouse service and user management

- Strong security and privacy with SOC 2 Type I compliance

**Serverless hosted ClickHouse for early adopters**

- Ecosystem of first-party connectors to onboard and work with data

- Automatic scaling up and down for compute

- Integrated billing using pay as you go pricing model

- Enhanced security features such as Private Link, IP Filtering, Auditing

**Serverless hosted ClickHouse for the broader market**

- Enhanced cloud console for analytics and operational controls

- Support for more AWS regions and AWS marketplace billing

- Uptime SLA and additional operations tools

- Advanced security features with SOC 2 Type II compliance

"

*Rokt has been an eager partner of ClickHouse as we modernize our analytics stack. By offloading operations to the experts our developers are focused on delivering the best experience possible while the business scales. We we are thrilled to see the path ClickHouse is forging.*

Attributor, role, company

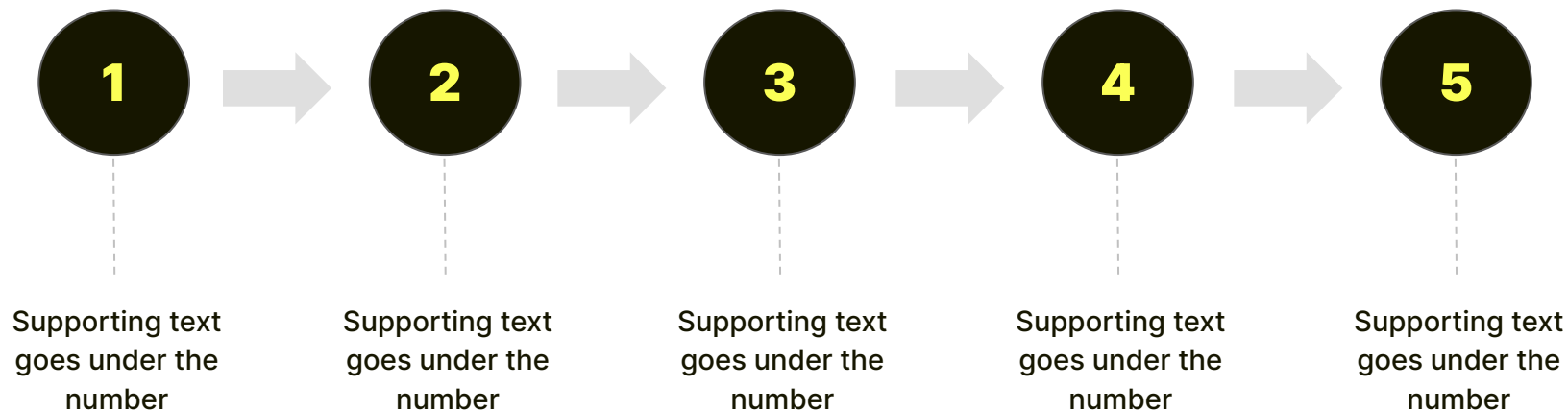ROKT

37

# Use this slide for code (light)

```
SELECT

    toStartOfMonth(upload_date) AS month,

    sum(view_count) AS `Youtube Views`,

    bar(sum(has_subtitles) / count(), 0.55, 0.7, 100) AS `% Subtitles`

FROM youtube

WHERE (month >= '2020-08-01') AND (month <= '2021-08-01')

GROUP BY month

ORDER BY month ASC


13 rows in set. Elapsed: 0.823 sec Processed 1.07 billion rows, 11.75 GB (1.30 billion rows/s., 14.27 GB/s.)
```
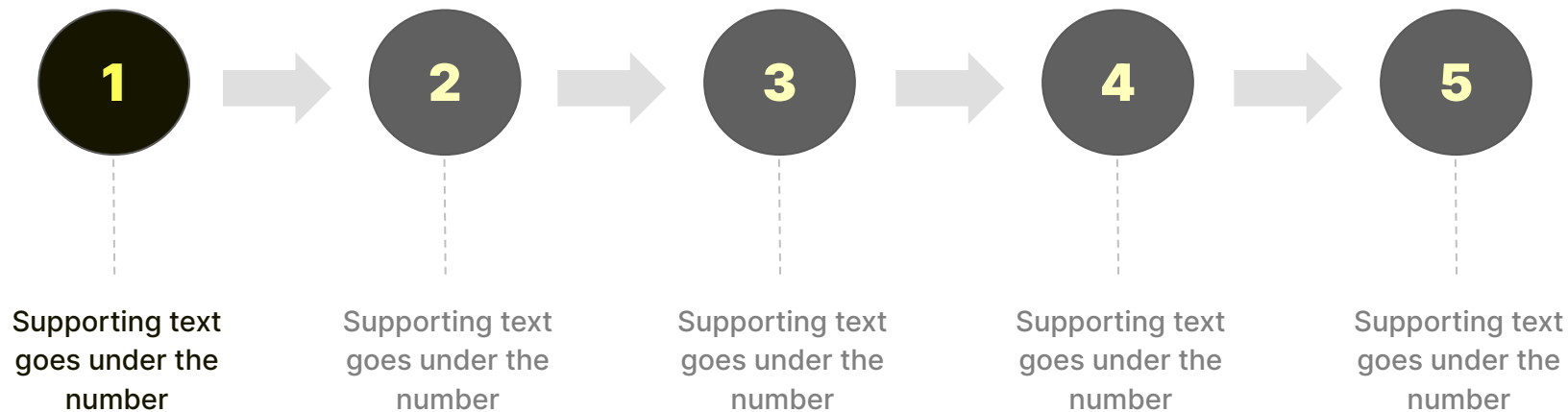
# Process diagram, 5 ideas

Here you could describe the topic of the section, or not

**1** → **2** → **3** → **4** → **5**

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

# Process diagram, 5 ideas

Highlighting one of the steps

**1** → **2** → **3** → **4** → **5**

**Supporting text goes under the number**

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

# Big number treatment

## 2.3k

### Header here

Supporting text goes here, under the header

## 1.2M

### Header here

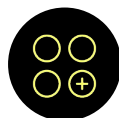Supporting text goes here, under the header

## 45

### Header here

Supporting text goes here, under the header

# Features

## Blazing fast

Uses all available hardware to its full potential to process each query as fast as possible. Peak processing performance for a single query stands at more than 2 terabytes per second.

## Fault tolerant

Supports async replication and can be deployed across multiple datacenters. All nodes are equal, which allows avoiding having single points of failure.

## Easy to use

ClickHouse is simple and works out-of-the-box. Simplifies data processing by instantly processing structured data using a user-friendly SQL dialect and eliminating non-standard API requirements.
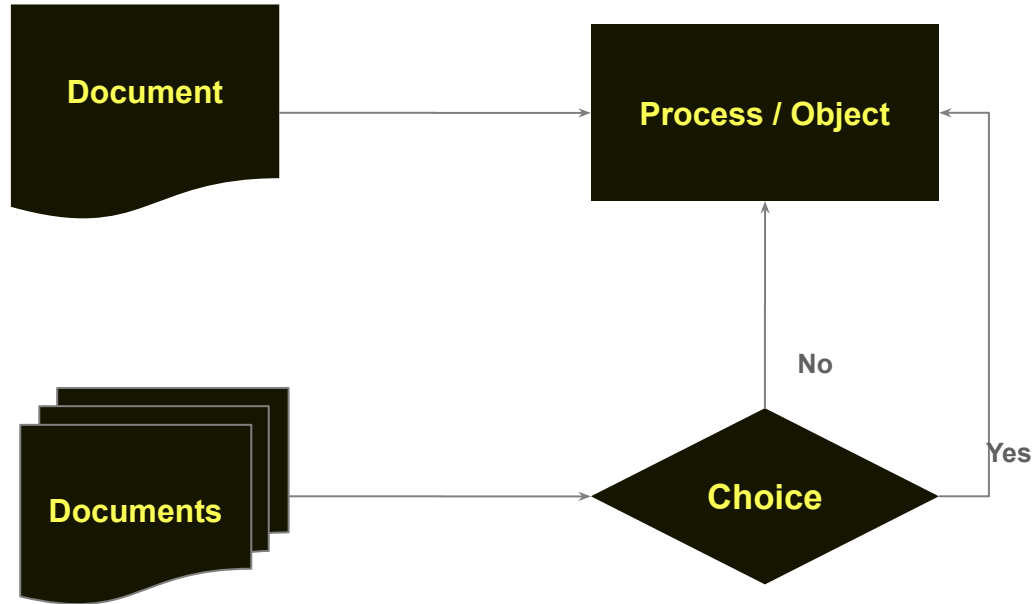
## Highly reliable

Can be configured as a purely distributed system located on independent nodes, without any single points of failure. It also includes a lot of enterprise-grade security features and fail-safe mechanisms against human errors.
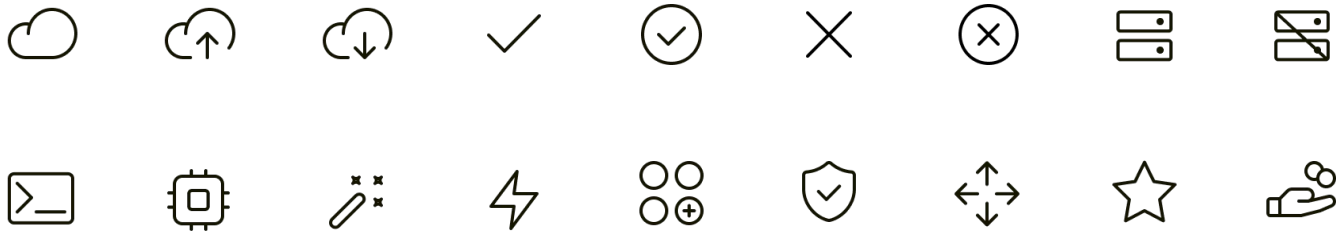
# Table format

| Header 1 | Header 2 | Header 3 | Header 4 | Header 5 |
|----------|----------|----------|----------|----------|
|          |          |          |          |          |
|          |          |          |          |          |
|          |          |          |          |          |
|          |          |          |          |          |
|          |          |          |          |          |
|          |          |          |          |          |

# For architecture/business process drawings

# Icons for light background

45

# Speakers (dark mode)

**Aaron Katz**

CEO @ ClickHouse

**Alexey Milovidov**

CTO @ClickHouse

**Yury Izrailevsky**

President and
VP of Engineering

# Table of contents - many topics (dark mode)

**01** **Content title**
Presenter / brief description

**02** **Content title**
Presenter / brief description

**03** **Content title**
Presenter / brief description

**04** **Content title**
Presenter / brief description

**05** **Content title**
Presenter / brief description

**06** **Content title**
Presenter / brief description

# Table of contents - fewer topics(dark mode)

**01** **Content title**
Presenter / brief content description

**02** **Content title**
Presenter / brief content description

**03** **Content title**
Presenter / brief content description

**04** **Content title**
Presenter / brief content description

# Section title

**Subtitle or description of the section.
Remove if not needed**

ClickHouse

# 01 Section title

**Subtitle or description of the section. Remove if not needed**

ClickHouse

# Content heavy slide

# Content slide - less text heavy

# Feature highlight

- Item 1
- Item 2
- Item 3



**another one** Running

| monitoring-internal | Directly assigned permissions |
| operator-internal | Directly assigned permissions |
| sql-console | default_role |
| default | default_role |

**My service** Running

| monitoring-internal | Directly assigned permissions |
| operator-internal | Directly assigned permissions |
| default | default_role |
| sql-console | default_role |
| sql-console:cristina.albu@clickhouse.com | default_role |

# Two cards highlight

## Title here

**Subtitle or a brief description**

- Another bullet here
- Another bullet here
- Another bullet here
- Another bullet here

## Title here

**Subtitle or a brief description**

- Another bullet here
- Another bullet here
- Another bullet here
- Another bullet here

# Three cards highlight

## Title here

**Subtitle or a brief description**

- Another bullet here
- Another bullet here
- Another bullet here
- Another bullet her

## Title here

**Subtitle or a brief description**

- Another bullet here
- Another bullet here
- Another bullet here
- Another bullet here

## Title here

**Subtitle or a brief description**

- Another bullet here
- Another bullet here
- Another bullet here
- Another bullet her

"

*Rokt has been an eager partner of ClickHouse as we modernize our analytics stack. By offloading operations to the experts our developers are focused on delivering the best experience possible while the business scales. We we are thrilled to see the path ClickHouse is forging.*

Attributor, role, company

ROKT

# Features

## Blazing fast

Uses all available hardware to its full potential to process each query as fast as possible. Peak processing performance for a single query stands at more than 2 terabytes per second.

## Fault tolerant

Supports async replication and can be deployed across multiple datacenters. All nodes are equal, which allows avoiding having single points of failure.

## Easy to use

ClickHouse is simple and works out-of-the-box. Simplifies data processing by instantly processing structured data using a user-friendly SQL dialect and eliminating non-standard API requirements.
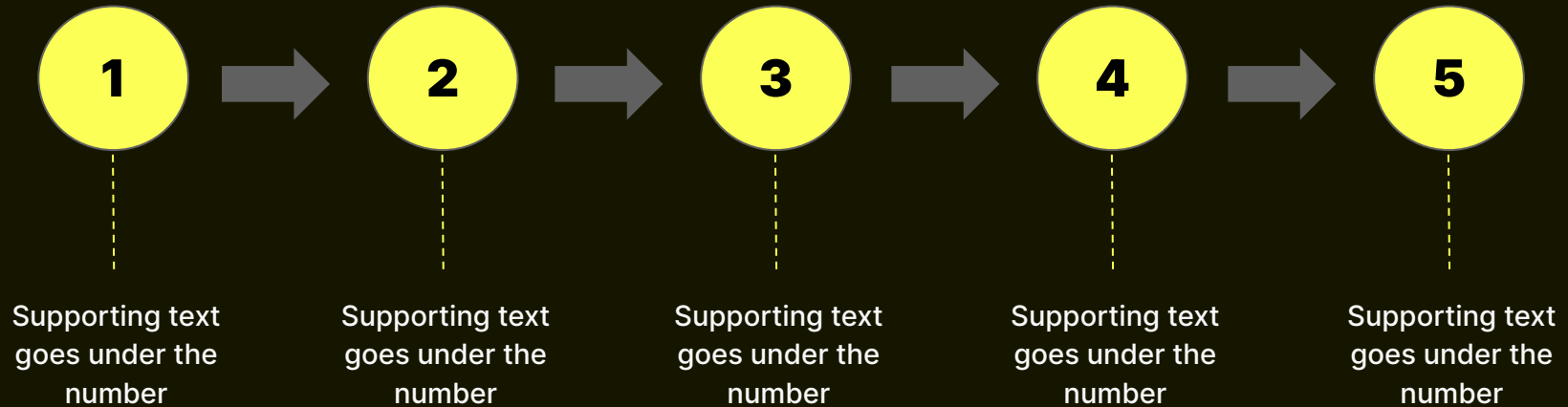
## Highly reliable

Can be configured as a purely distributed system located on independent nodes, without any single points of failure. It also includes a lot of enterprise-grade security features and fail-safe mechanisms against human errors.
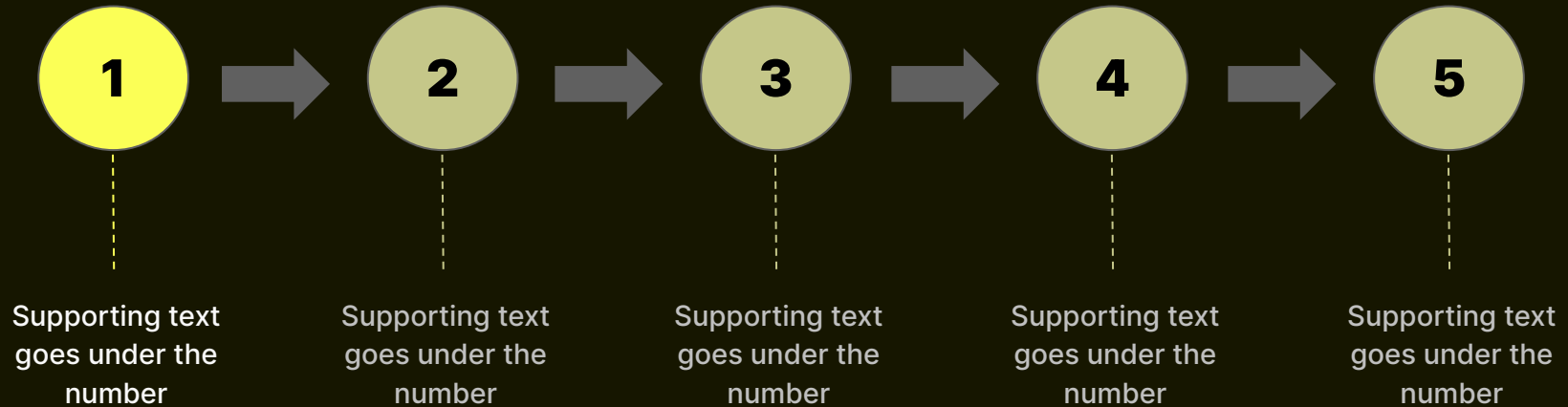
# Process diagram, 5 ideas

Here you could describe the topic of the section, or not
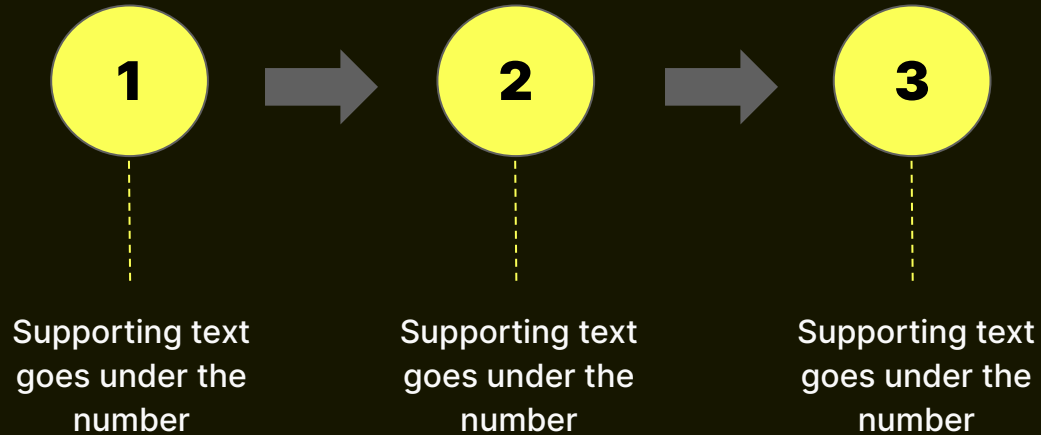
| 1 | → | 2 | → | 3 | → | 4 | → | 5 |

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

# Process diagram, 5 ideas

Highlighting one of the steps

**1** → **2** → **3** → **4** → **5**

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

# Process diagram, 3 ideas

Here you could describe the topic of the section, or not

**1**  →  **2**  →  **3**

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

# Process diagram, 3 ideas

Highlighting one of the ideas

**1** → **2** → **3**

Supporting text goes under the number

Supporting text goes under the number

Supporting text goes under the number

# Big number treatment

## 1.2k

### Header here

Supporting text goes here, under the header

## 1M

### Header here

Supporting text goes here, under the header

## 45

### Header here

Supporting text goes here, under the header
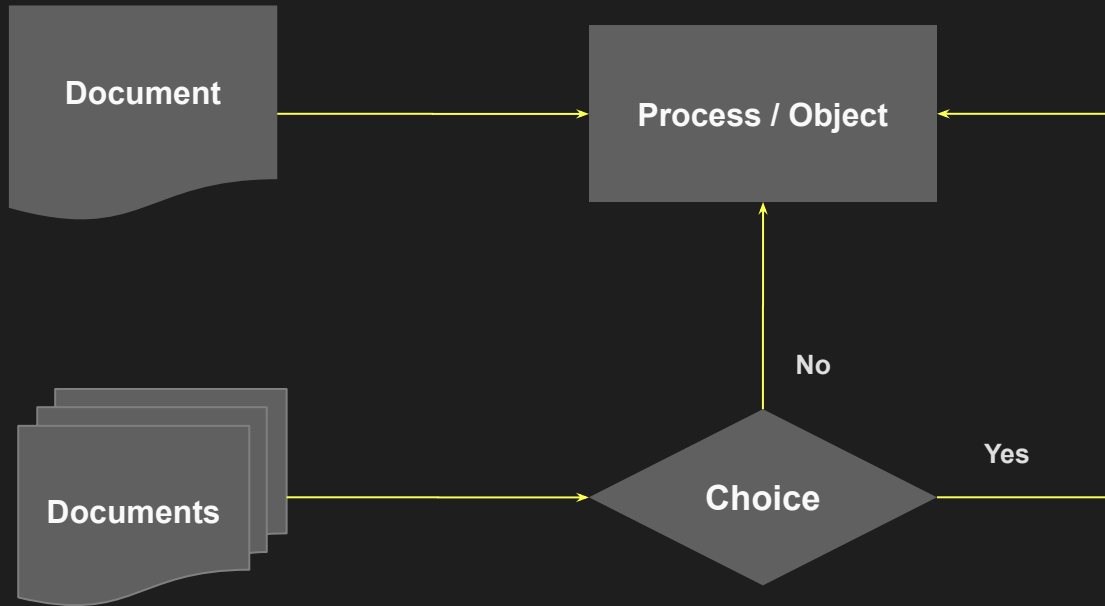
# Use this slide for code

Use template colors to highlight code. Feel free to remove this paragraph.

```
SELECT

    toStartOfMonth(upload_date) AS month,

    sum(view_count) AS `Youtube Views`,

    bar(sum(has_subtitles) / count(), 0.55, 0.7, 100) AS `% Subtitles`

FROM youtube

WHERE (month >= '2020-08-01') AND (month <= '2021-08-01')

GROUP BY month

ORDER BY month ASC


13 rows in set. Elapsed: 0.823 sec Processed 1.07 billion rows, 11.75 GB (1.30 billion rows/s., 14.27 GB/s.)
```

# For architecture / business process drawings

# Icons for dark background