

# Taming the Data Beast Using DataHub

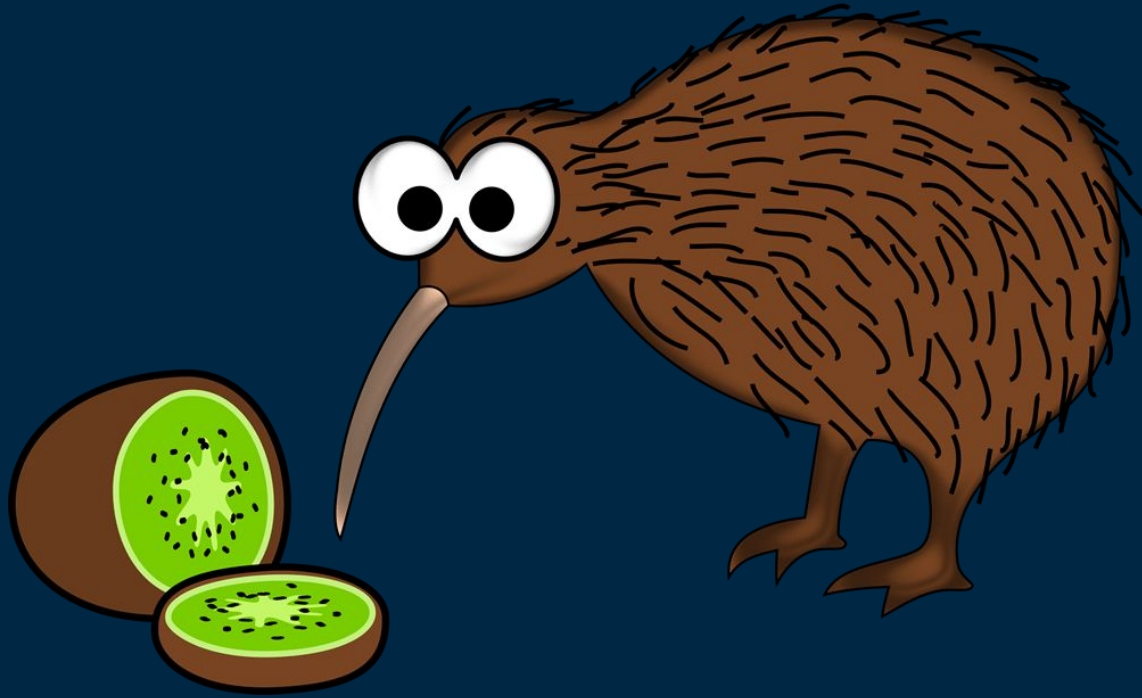
Data Engineering Melbourne Meetup  
Nov 25, 2020  
Mars Lan

# MARS LAN



- Co-creator of DataHub
- CTO @ Metaphor Data - 1w
- TL @ LinkedIn Metadata Team - 4y
- SWE @ Google (GCP, Android) - 3.5y
- UCLA CS PhD
- Twitter: @mars\_lan

I'M A KIWI!



The background is a dark blue gradient. It is decorated with an abstract pattern of small squares and thin vertical lines in various colors including teal, orange, pink, and light blue. These elements are scattered across the slide, with some appearing as solid shapes and others as outlines. The text is centered in the middle of the slide.

“Software is eating the world”

—Marc Andreessen, 2017

The background is a dark blue field decorated with a sparse, abstract pattern of geometric elements. These include small squares in various colors (light blue, orange, pink, teal) and thin white vertical lines of varying lengths, scattered across the frame.

“Data is eating the world”

—Everyone, 2020

# MACRO TRENDS

- Data Democratization ⇒ More Organic
  - Data mesh, decentralized data governance, self-service, remote work
- Role Specialization ⇒ More Personas
  - Data scientists/analysts/engineers, AI/ML engineers, business analysts/users, ...
- Explosion of Data Systems & Tools ⇒ More Complexity
  - Hadoop, Spark, Flink, Kafka, Presto, TensorFlow, Elasticsearch, MongoDB
- Adoption of Cloud Computing ⇒ More Data
  - Easier & cheaper than ever to create more data
- Increasing Regulatory Pressure ⇒ More Controls
  - GDPR, CCPA, LGPD, BCBS 239, MiFID II, FRTB, CCAR

# PROBLEMS

- Finding data is hard
  - Data lake  $\Rightarrow$  Data swamp
  - Siloed teams  $\Rightarrow$  Siloed data
  - Specialized systems  $\Rightarrow$  Specialized data
- Managing data is hard
  - Governance
  - Compliance
  - Policy-driven data management
- Trusting data is hard
  - Lineage
  - Data availability
  - Data quality & health
  - Data profiling & distribution
  - Ownership & documentation
  - Certification & curation

The background is a dark blue gradient. It is decorated with various geometric elements: thin white vertical lines of varying lengths, small squares in teal, orange, and pink, and larger squares in teal and orange. The text is centered in a light orange, sans-serif font.

Metadata,  
Metadata,  
Metadata



# WHAT IS METADATA?

**Who** created this?

**When** was it last updated?

**What** does each column mean?

|   | A         | B           | C           | D            |  |
|---|-----------|-------------|-------------|--------------|--|
| 1 | <b>ID</b> | <b>Name</b> | <b>Date</b> | <b>Value</b> |  |
| 2 | 7792      | June        | 2013/05/14  | 4            |  |
| 3 | 2675      | April       | 2020/09/01  | 0            |  |
| 4 | 4190      | Joe         | 1987/12/2   | NULL         |  |
| 5 | 3655      | May         | 2005/11/17  | 3            |  |
| 6 | ...       | ...         | ...         | ...          |  |

**Where** did data come from?

**Why** is there NULL value?

**How** was Value column computed?

# LINKEDIN METADATA JOURNEY

Pull-based (crawlers)  
Monolithic app  
Table-based models

OSS WhereHows (V1)

2016

2017/18

WhereHows + TMS (V2)

Push-based (Kafka)  
App + monolithic service  
Opinionated models

Push-based (Kafka)  
App + distributed services  
Generalized models

DataHub + GMA (V3)

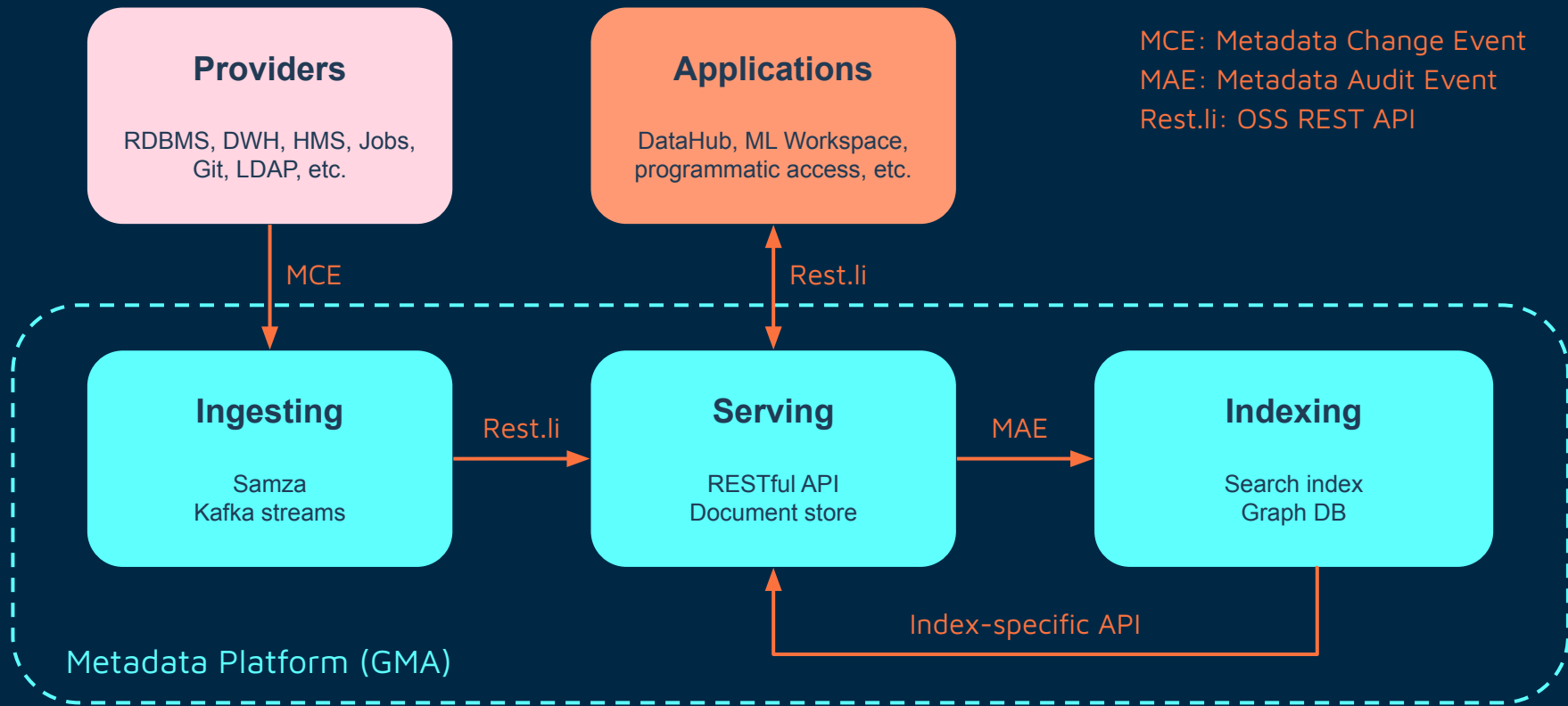
2019

2020

OSS DataHub

Docker & K8s  
Cloud integration  
GraphQL

# DATAHUB ARCHITECTURE



# METADATA PLATFORM (GMA)

- Scalable

- Web-scale stack
- Distributed storage, indexing & serving
- four/five 9's uptime
- Decentralized metadata modeling

- Enrichable

- Beyond read-only aggregator
- Collaborative edits
- Human curation

- Queryable

- Key-value
- Distributed joins
- Full-text search
- Graph traversal

- Real-time

- Stream-based ingestion
- Event-driven architecture
- Trigger-based applications

# WHY STREAMS (KAFKA)?

- Near real-time
  - $O(\text{seconds})$  delay
- Loose coupling
  - Non-blocking, fire-and-forget
- Queuing
  - Smooth out bursty traffic
  - Async consumption
- Schema Compatibility
  - Easy to enforce via schema registry
- Scalable
  - Multi-readers/writers
  - Partition-level parallelism
- Persistent storage
  - Sequential key-value store
  - Bootstrap & backfill

# INGESTION MODES

- Existing metadata services

- Crawler: uninstrumentable
- Direct event emission: instrumentable
- Event conversion: existing events

- New metadata services

- DAO: “man-in-the-middle” integration

- Metadata in Git

- Build-time: tooling emits event
- Publish-time: events artifact
- Deploy-time: services/jobs emit events

# METADATA MODELING

## Nodes

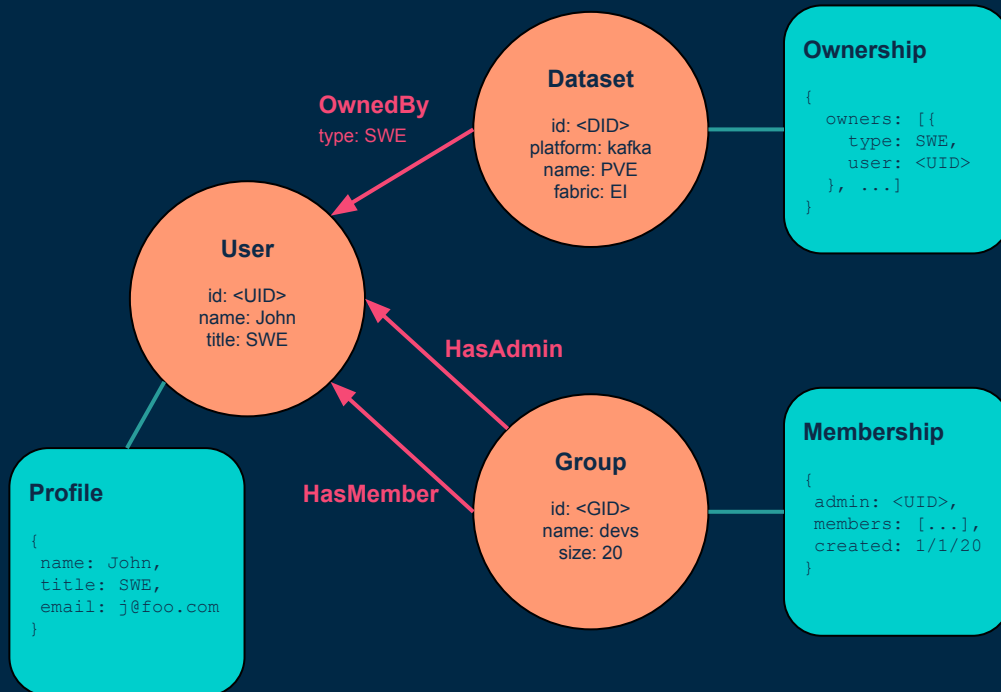
Entities, e.g. datasets, partitions, features, users, groups, experiments, ...

## Edges

Relationship, e.g. OwnedBy, DerivedFrom, Contains, HasMember, ...

## Documents

Metadata, e.g. ownership, membership, upstreams, configs, compliance metadata, ...



# ADOPTION

- LinkedIn (within 18 months)

- Integrated with 40 teams/projects
- 30+ entities, 200+ types of metadata
- Use cases
  - Search & discovery
  - Data privacy compliance
  - Access control
  - Life-cycle management
  - Data Ops
  - AI DevOps

- OSS (within 6 months)

- 8 companies running in production
- 20 companies building POC or seriously evaluating
- Success stories
  - Expedia: Data-driven tech company
  - Saxo Bank: Investment bank going through digital transformation
  - SpotHero: Cloud-native small startup (3 data engs & 50 data users)



# OSS ROADMAP

- New Entities & Relationships

- Jobs & flows\*
- Dashboards\*
- AI features/models\*
- Business glossary\*
- Schemas
- Metrics
- Services

- Integrations

- BI: Looker\*, Mode, Redash, Superset
- Scheduler: AirFlow, Dagster, Azkaban
- Data Quality: Great Expectations

- Features

- Entity insights
- Data privacy
- Governance
- SSO: OIDC & SAML

- Platform

- Gremlin-based query DAO
- Aspect-specific events
- GraphQL API
- NoSQL backend (e.g. MongoDB)
- OLAP index (e.g. Pinot, Druid)



## Data concepts related to "foodie"

[Browse All Data Concepts](#)

Curated information about key data terms presented with relevant context of known data entities.

## DATA CONCEPT

**Engaged Quality Foodies**

A Quality Foodie (QF) is a member who satisfies the criteria of: Reachable - Can be contacted by restaurants about the promotions or new menu item...

## Filters

## Data Origin ⓘ

☐ Prod (15)☐ Corp (12)

## Platform

☐ Mysql (24)☐ Hdfs (3)

## Datasets

Showing 1 - 10 of 27 results

[/demo/teamY/Foodie](#)

Data Origin PROD  
Platform hdfs  
Health **100%**

0 - 0 ⓘ

[/demo/teamX/Foodie](#)

Data Origin PROD  
Platform hdfs  
Health **100%**

0 - 0 ⓘ

[/nacho/test/Foodie](#)

Data Origin PROD  
Platform hdfs  
Health **100%**

0 - 0 ⓘ



## Datasets

/demo/rolling\_aggregate\_orders\_2 Hdfs UndefinedFabric: **HOLDEM/WAR**

## Health

Last calculated a month ago

100%

[See Details](#)

0 - 0

Schema

Status

ACL Access

**Ownership**

Compliance

Dataset Groups

Relationships

Health

Docs

Last Saved: 3 months ago by nkanamar

## Owners

Please maintain at least 2 owners.

| LDAP Username                  | Full Name               | ID Type | Ownership Type |  |
|--------------------------------|-------------------------|---------|----------------|--|
| nkanamar                       | Nagarjuna Kanamarlapudi | USER    | DataOwner      |  |
| pgunnam                        | Pardhu Gunnam           | USER    | DataOwner      |  |
| <a href="#">+ Add an owner</a> |                         |         |                |  |

Please make changes to save.

Save

[Metrics](#) > [all](#) > [foodie](#) > [app\\_funnel](#) > [activated](#)

## METRICS

**activated**

number of users activated the app

## Owners

[malan](#), [pgunnam](#)Health <sup>®</sup>

Last calculated 5 months ago

100%

[See Details](#)[Overview](#)[Related Entities](#)[Health](#)[Docs](#)

## Related entities

This metric is derived from [1 dataset](#) and related to [17 metrics](#)[View the sections below for more details](#)

## Derived from datasets (1)

| Name                       | Description             | Owners  |
|----------------------------|-------------------------|---|
| <a href="#">app_funnel</a> | All lite funnel metrics | <a href="#">pawkumar</a> , <a href="#">jmodi</a> , <a href="#">adchoudh</a> , <a href="#">fmt</a> , <a href="#">kisir</a> |

## Other metrics from the same dataset (10)

| Name                                   | Alias | Description                              |
|--|-------|--|
| <a href="#">download</a>               | -     | number of users downloaded the app       |
| <a href="#">signup_pass_overall</a>    | -     | number of users signup overall           |
| <a href="#">first_time_signup_pass</a> | -     | number of users signup in the first time |
| <a href="#">first_time_signin</a>      | -     | number of users first time logged in     |
| <a href="#">first_signin_failed</a>    | -     | number of users first login failed       |



## Mars Lan

Staff Software Engineer, Systems Infrastructure

[Edit profile](#)

### ASK ME ABOUT

[metadata](#)

TL for the metadata team

### TEAM

[datahub](#)[gma](#)

### MANAGER



Tai Tran

[malan@linkedin.com](mailto:malan@linkedin.com)[LinkedIn Profile](#)[Cinco Profile](#)

#### Access Management

JIT Datasets

#### Lists

Features

Liked Entities

Followed Entities

#### Data you might be interested in

Datasets

#### Ownership

Datasets

#### Metrics

Charts

Dashboards

ML features

UMP Flows

## Metrics Mars Lan Owns

Showing 1 - 10 of 17 results

### count\_of\_events

number of events produced

|           |   |
|-----------|---|
| Bucket    | datahub                                 |
| Formula   | COUNT                                   |
| Dataset   | datahub_page_view_event                 |
| XLNT Tier | Daily: 2                                |
| Tags      | datahub, metadata, pageviewevent, route |
| Frequency | DAILY                                   |
| Health    | 100%                                    |

### count\_of\_events

number of events produced

|           |  |
|-----------|--|
| Bucket    | datahub                                |
| Formula   | COUNT                                  |
| Dataset   | datahub_search_impression_event        |
| XLNT Tier | Daily: 2                               |
| Tags      | datahub, metadata, search, impressions |
| Frequency | DAILY                                  |
| Health    | 100%                                   |

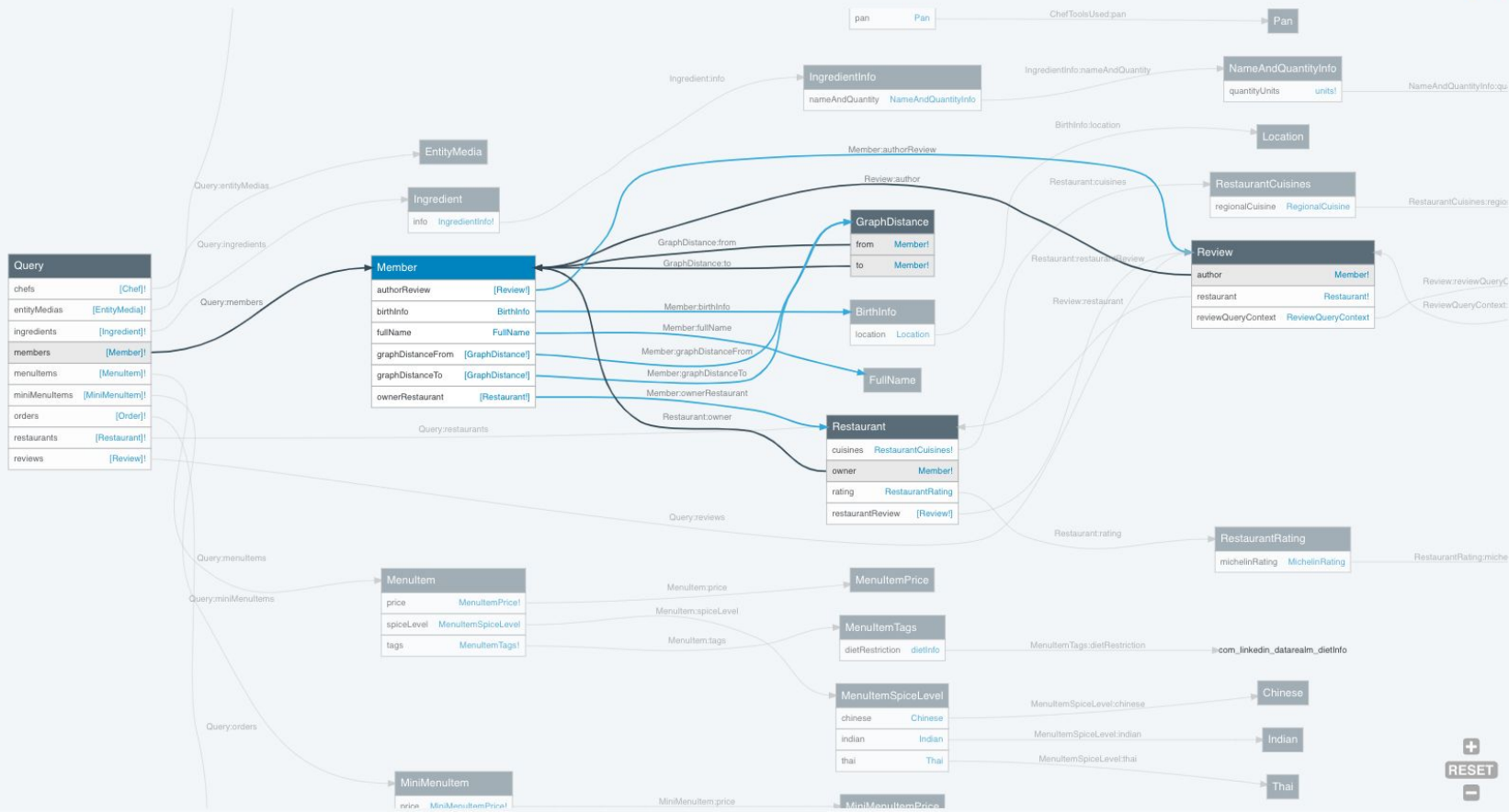
### mie\_total\_count

number of mie events

|        |     |
|--------|-----|
| Bucket | tms |
|--------|-----|



## Schema Graph View

Schemas [datarealm-restaurants.opencrud.graphql](#)

| /demo/rolling_aggregate_orders_2    |
|-------------------------------------|
| /city                               |
| /date                               |
| /restaurantName                     |
| /no_of_orders                       |
| <a href="#">View dataset detail</a> |

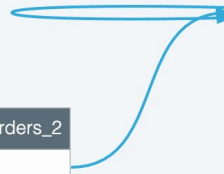
| /demo/aggregate_restaurant_orders_2 |
|-------------------------------------|
| /city                               |
| /date                               |
| /restaurantName                     |
| /no_of_orders                       |
| <a href="#">View dataset detail</a> |

| /demo/dim_customer_2                |
|-------------------------------------|
| /address/city                       |
| /address/city                       |
| <a href="#">View dataset detail</a> |

| demo.customer_info_changes_2        |
|-------------------------------------|
| /city                               |
| <a href="#">View dataset detail</a> |

| /demo/fact_orders_2                 |
|-------------------------------------|
| /timeStamp                          |
| /restaurantName                     |
| /orderNumber                        |
| <a href="#">View dataset detail</a> |

| demo.kafka_orders_2                 |
|-------------------------------------|
| /header/timeStamp                   |
| /restaurantName                     |
| /orderNumber                        |
| <a href="#">View dataset detail</a> |



# GitHub: linkedin/datahub

Monthly town hall

Slack workspace

PRs welcome!





Do you have any questions?

[mars@metaphor.io](mailto:mars@metaphor.io)  
[linkedin.com/in/marslan](https://linkedin.com/in/marslan)  
Twitter: [mars\\_lan@](#)

# THANKS

DISCLAIMER: A large portion of this deck is based on the published Budapest Data Forum 2020 talk I gave as an employee of LinkedIn.

CREDITS: This presentation template was created by [Slidesgo](#),  
including icons by [Flaticon](#), and infographics & images by [Freepik](#)  
Please keep this slide for attribution