# MICROSOFT FABRIC

**Bas Land**

# Your first Lakehouse using Microsoft
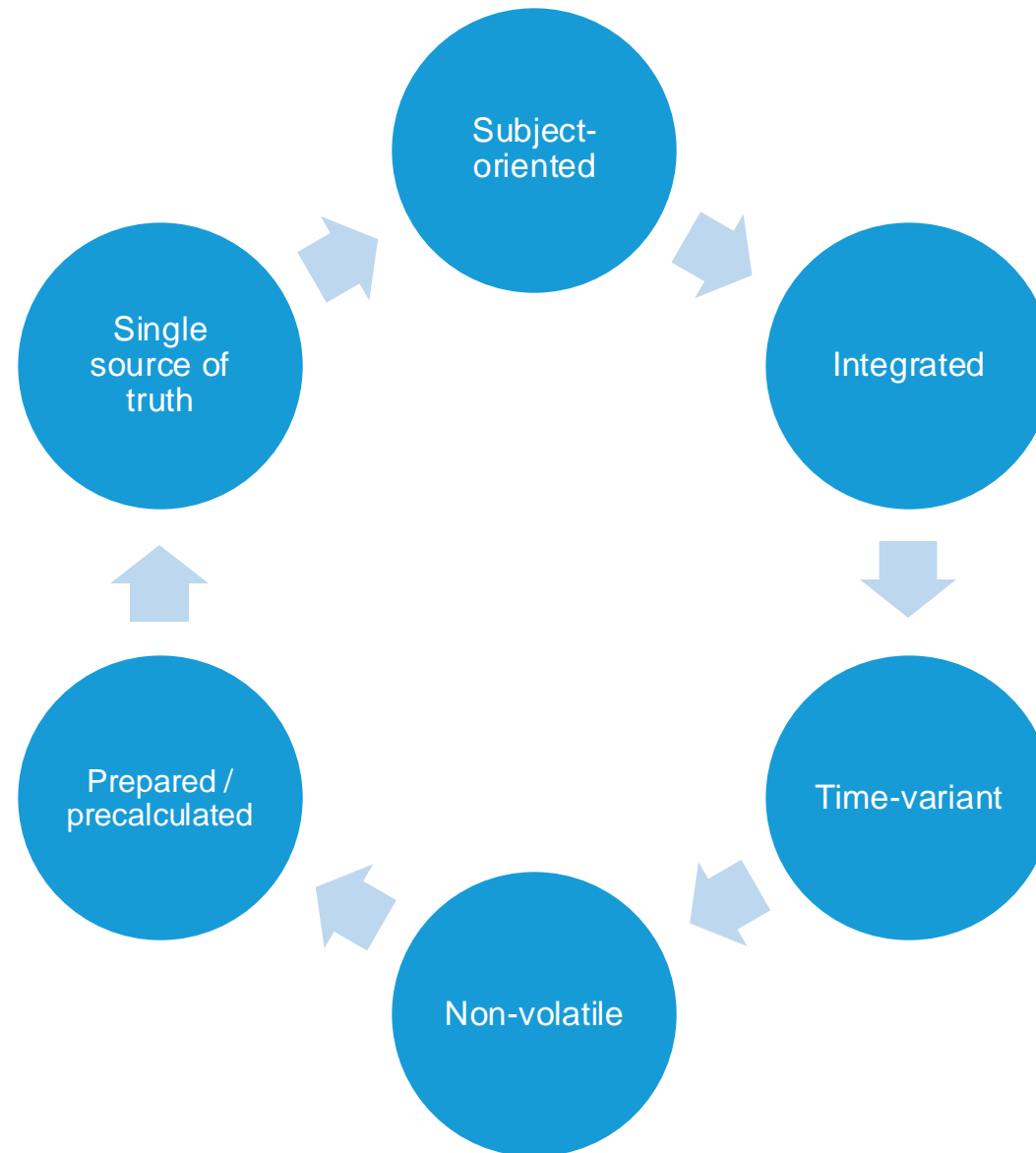
# Today

- Core concepts for data warehousing in Fabric

- Your first Fabric Lakehouse:
  - Rulebook
  - ELT process
  - End-to-end solution design with examples

- Takeaways

# Core concepts

# General design principles

# Lakehouse: under the hood

| Ingest | Store | Process |
|--------|-------|---------|
| Pipelines<br>Notebooks<br>Dataflows | OneLake<br>Data lake blob storage | Notebooks<br>(PySpark) |

# Fabric concepts

1. Tenant design

2. Storage

3. Compute

# Fabric: tenant design

Tenant & OneLake

Domains

Workspaces

Folders

Artefacts

# Fabric: tenant design
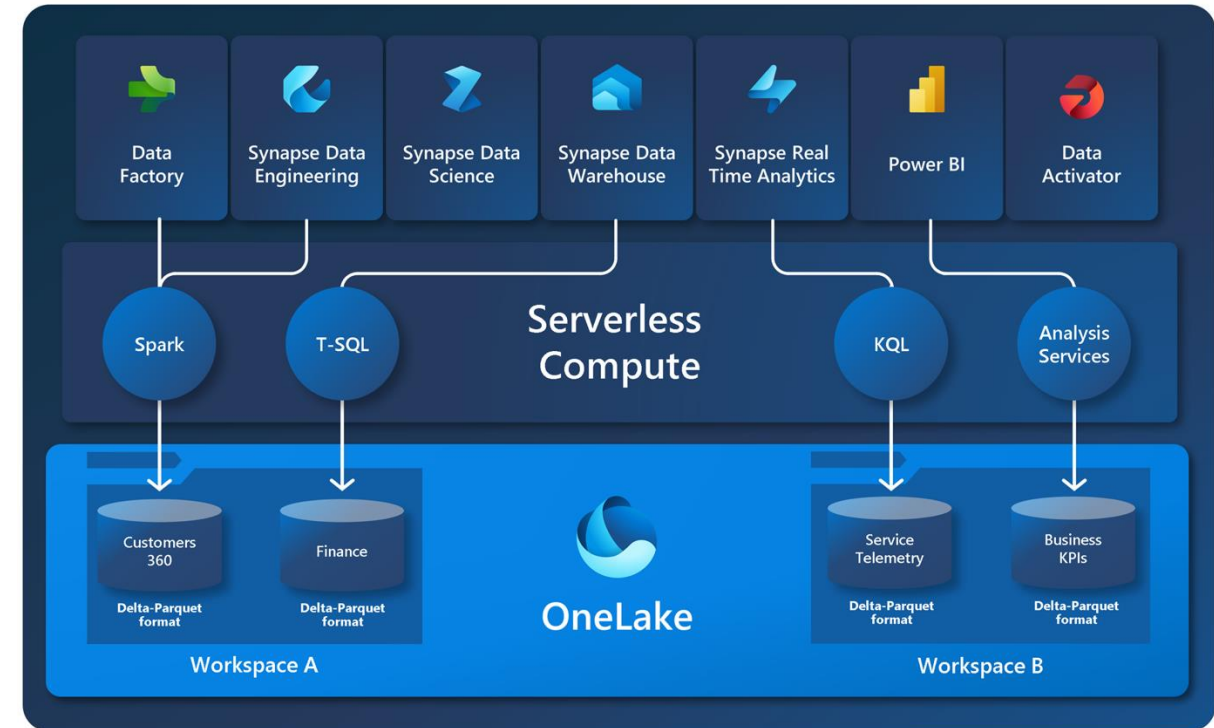
# Fabric: storage

- Storage in Fabric is called 'OneLake'
- It's data lake storage (blob), with one OneLake per organisation (tenant)
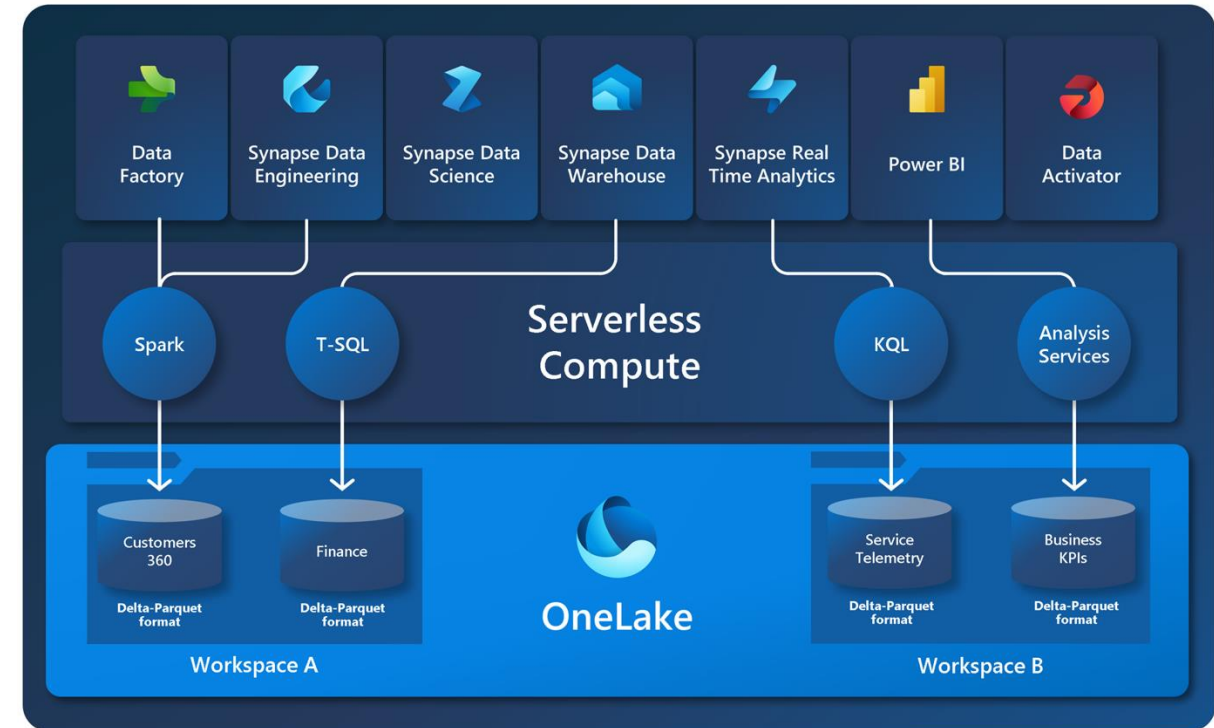- Storage is cheap! ± 2,2 cent per GB per month

# Fabric: compute

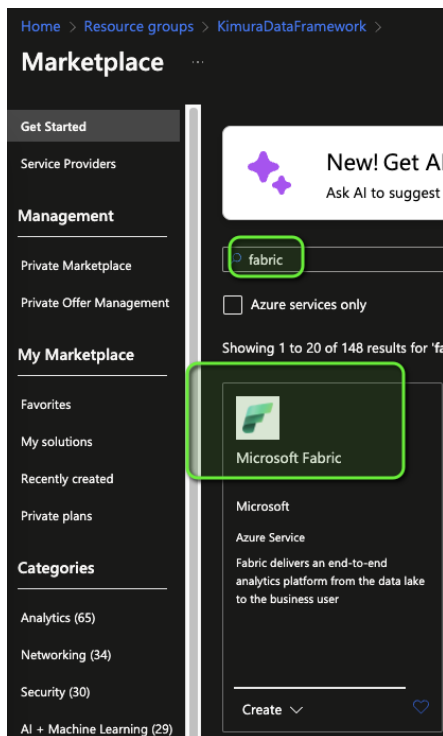- Compute in Fabric is Serverless, and included in your SaaS capacity

- For data engineering purposes, you get a Spark cluster that is preconfigured and can be configured to your custom needs

- Compute = expensive. Be smart about it ☺

# Buying the Fabric compute capacity

- portal.azure.com
- I'd create at least a new Resource Group, maybe even a Subscription (depending on organisation requirements)

# Buying the Fabric compute capacity

- After deployment, configure the Capacity Administrator:
- Unfortunately, no Entra ID groups…

# Creating the first Fabric Workspace

- Head over to app.powerbi.com

- Log in with the same account that is the Capacity Administrator

- Create a new workspace,

- Now, pay attention to the License Type

# Creating the first Fabric Workspace

# Medallion architecture

# Workspaces and Medallion

- One workspace per Medallion stage

- One workspace per DTAP environment

# Your first Fabric Lakehouse

# The scenario



- Client: small B2B professional services company
- Main software: Moneybird, Odoo
- Data domains: Finance, Projects & Timesheets

# The rulebook

- Data Factory Pipelines orchestration

- PySpark notebooks for processing

- Reusable code for everything

- Lakehouse for storage

# Prerequisites

Setting up:

- Fabric compute capacity (done)

- Fabric workspaces (done)

- Configuring Spark (we skip that for today)

- Setting up the Environment (**important**!)

# Environments in Fabric

- Fabric contains a lot of Python packages natively

- You can also add packages using the PyPI (Python Package Index)

- And, you can write and include your own code

- This is done using Environments

# Environments in Fabric

- Create a new environment:

# Environments in Fabric

- Add code from PyPI:

# Environments in Fabric

- Add your own code:

# Three lakehouses



- Bronze: timestamped folders, (incremental loads), raw data (files; json, xml, csv, parquet etc)

- Silver: historical archive, normalised schemas, delta tables

- Gold: prepared and modeled DWH tables, dims and facts, delta tables

# Setting up the ELT process



- Keep in mind the Medallion architecture

- Start connecting to data, to extract Moneybird and Odoo to Bronze

- We do so using Python, executed in Fabric Notebooks, scheduled and run in Fabric Pipelines

# ELT: Extract

**Data Factory Pipelines**
- Orchestrate entire process

**Orchestrator notebook**
- Execute worker notebook
- runMultiple() for parallel executions

**Worker notebook**
- Execute generic function

**Reusable Python functions in the Kimura Data Framework (KDF)**
- Connect to source API
- Save JSON response to the lakehouse

# ELT: Bronze folder structure

Explorer «

∨ Bronze_lakehouse

  ∨ 🗁 Tables

  ∨ 🗁 Files

    ∨ 🗁 Landing

      ∨ 🗁 Odoo

        ∨ 📁 202409031641

          › 🗁 account_analytic_line

          › 🗁 hr_employee

          › 🗁 project_project

          › 🗁 project_task

          › 🗁 res_partner

          › 🗁 res_users

- Root folder 'Landing' (for shortcuts!)

- Source system

- Timestamp

- Folder per table/endpoint

- Files managed by Fabric!

# Silver lakehouse

## Explorer

- **Silver_lakehouse**
  - Tables
    - silver_0_ExtractLogging
    - silver_account_analytic_line
  - Files
    - Landing

## silver_account_analytic_line

| | ABC account_id | 0/1 allow_billable | 1.2 amount | 0/1 auto_accoun... | ABC category | 0/1 code |
|---|---|---|---|---|---|---|
| 1 | [1,"Internal"] | False | -505.76 | False | other | False |
| 2 | [1,"Internal"] | False | 0 | False | other | False |
| 3 | [1,"Internal"] | False | -505.76 | False | other | False |
| 4 | [1,"Internal"] | False | 0 | False | other | False |
| 5 | [1,"Internal"] | False | -505.76 | False | other | False |
| 6 | [1,"Internal"] | False | 0 | False | other | False |
| 7 | [1,"Internal"] | False | 0 | False | other | False |
| 8 | [1,"Internal"] | False | 0 | False | other | False |
| 9 | [1,"Internal"] | False | 0 | False | other | False |
| 10 | [1,"Internal"] | False | 0 | False | other | False |
| 11 | [1,"Internal"] | False | 0 | False | other | False |
| 12 | [1,"Internal"] | False | 0 | False | other | False |
| 13 | [1,"Internal"] | False | 0 | False | other | False |
| 14 | [159,"V/2024/00... | True | -18.42 | False | other | False |
| 15 | [149,"V/2024/00... | True | -36.84 | False | other | False |
| 16 | [164,"Internal"] | False | -189.66 | False | other | False |
| 17 | [1,"Internal"] | False | -71.88 | False | other | False |
| 18 | [142,"Sales & Bu... | False | -63.22 | False | other | False |
| 19 | [149,"V/2024/00... | True | -64.47 | False | other | False |

# Gold lakehouse

Explorer

- ∨ Gold
  - ∨ 📁 Tables
    - › ⊞ dimCustomers
    - › ⊞ dimDates
    - › ⊞ dimProjects
    - › ⊞ dimServices

## dimCustomers

| ⊞ | 123 CustomerId | ABC RelationType | ABC CustomerNa... | ABC IsActive | ABC PaymentType |
|---|---|---|---|---|---|
| 1 | 27 | Klant | | True | banktransfer |
| 2 | 37 | Klant af | | True | directdebit |
| 3 | 103 | Klant | | True | banktransfer |
| 4 | 104 | Klant | | True | banktransfer |
| 5 | 128 | Klant af | | True | directdebit |
| 6 | 130 | Klant | | True | banktransfer |

# End-to-end solution with examples

# End-to-end solution with examples

**Invoke Pipeline** ↗
- Process Bronze

**Set variable**
- (*x*) Set variable Datetime

**Copy data**
- Copy metadata

**Invoke Pipeline** ↗
- Extract Odoo
- Extract Odoo

**Invoke Pipeline** ↗
- Extract Moneybird
- Extract Moneybird

**Lookup**
- Lookup Metadata

**Notebook**
- Bronze orchestrator

# End-to-end solution with examples

Lookup
```
Lookup Metadata
```
→
Notebook
```
Bronze orchestrator
```

```
1   if output_filename == 'financial_mutations':
2       mb.get_financial_mutations(
3           api_endpoints=api_endpoints.split(','),
4           api_tokens=api_tokens.split(','),
5           landing_path=landingzone_path,
6           output_filename=output_filename
7       )
8   else:
9       mb.get_endpoint(
10          api_endpoints=api_endpoints.split(','),
11          api_tokens=api_tokens.split(','),
12          landing_path=landingzone_path,
13          output_filename=output_filename
14      )
```
- Command executed in 24 sec 371 ms by Kevin Land on 11:12:50 AM, 8/26/24

```
3
4   for index,
5       endpoi
6
7       for adr
8           en
9
10          la
11
12          ext
13
14          if
15
16          if
17
18
19
20
21          en
22
23      api_en
24
25      activi
26          'n
27          'p
28          't
29          'a
30
31
32
33
34          }
35      }
36
37      activi
38
39  DAG = {
40      "activ
41  }
42
43  mssparkuti
44
45  new_timest
46  new_timest
47  new timest
```

```
api_endpoints = ','.join(endpoints)

    activity = {
        'name': row['output_filename'],
        'path': 'MB extract GET',
        'timeoutPerCellInSeconds': 9000,
        'args': {
            'output_filename': row['output_filename'],
            'timestamp': timestamp,
            'api_endpoints': api_endpoints,
            'api_tokens': api_tokens
        }
    }

    activities.append(activity)

DAG = {
    "activities": activities
}

mssparkutils.notebook.runMultiple(DAG)

new timestamps = [Row(administration id=int(admin id), tir
```

# End-to-end solution with examples

```
1    #Mapping logic
2    gold_table = "dimServices"
3    business_key_columns = ['ServiceId']
4
5    df_mapping = spark.sql("""
6    select
7        t.`data.id` as ServiceId
8        ,t.`data.project_id` as ProjectId
9        ,t.`data.name` as ServiceName
10       ,t.`data.invoice_method` as InvoiceMethod
11       ,t.`data.amount` as Amount
12       ,t.`data.start_date` as StartDate
13       ,t.`data.end_date` as EndDate
14       ,t.`data.revenue_group.label` as RevenueGroup
15       ,t.`data.subscription_cycle` as InvoiceCycle
16       ,t.`data.price` as Price
17       ,t.`data.amount` * t.`data.price` as Value
18       ,case when t.`data.revenue_group.id` = 'revenuegroup:ba06e5bec6c1b6b5' then 'Licenses' else 'Other' end as MRRLabel
19       ,c.Sys_ID as FK_dimCustomers_Sys_ID
20   from Bronze.bronze_Odoo_ProjectsService t
21   left join Bronze.bronze_Odoo_ProjectsProject p
22       on t.`data.project_id` = p.`data.id`
23   left join Silver.silver_dimCustomers c
24       on p.`data.organization.id` = c.OdooId
25   where 1=1
26   """)
27
28   k.gold_load_table( \
29       gold_table = gold_table_name, \
30       df_mapping = df_mapping, \
31       business_key_columns = business_key_columns, \
32       spark_session = spark\
33   )
```

Recap

# Recap

1. Keep core-concepts in mind

2. Setup a rulebook and abide by it

3. Use tools and reusable code!

4. Fabric lakehouses aren't that scary ☺

# Bas Land

- BI consultant since 2013
- Co-founder of two data companies:
  - Kimura Data Intelligence (consultancy)
  - DataChimp (SaaS analytics for accounting firms)


- Married to Anouk, we have a daxhund (☺) called Chester
- Sports: purple belt Brazilian jiu-jitsu,
  weight lifting, running

# Kimura Data Intelligence B.V.

Fonteinkruid 6b
3931 WX, Woudenberg
The Netherlands

www.kimura.nl
info@kimura.nl

**Bas Land**
Managing Partner

bas@kimura.nl