



· Self-serve data ingestion for  
· high-scale workloads in the cloud

$\varphi$  upsolver

# Upsolver

## What we do

Upsolver is a self-serve cloud data ingestion service for high-scale workloads (*big data, streaming, AI*).

Product is purpose-built for the cloud with a Snowflake-like architecture for data processing.

## Certifications



## Partners



## Proven at scale

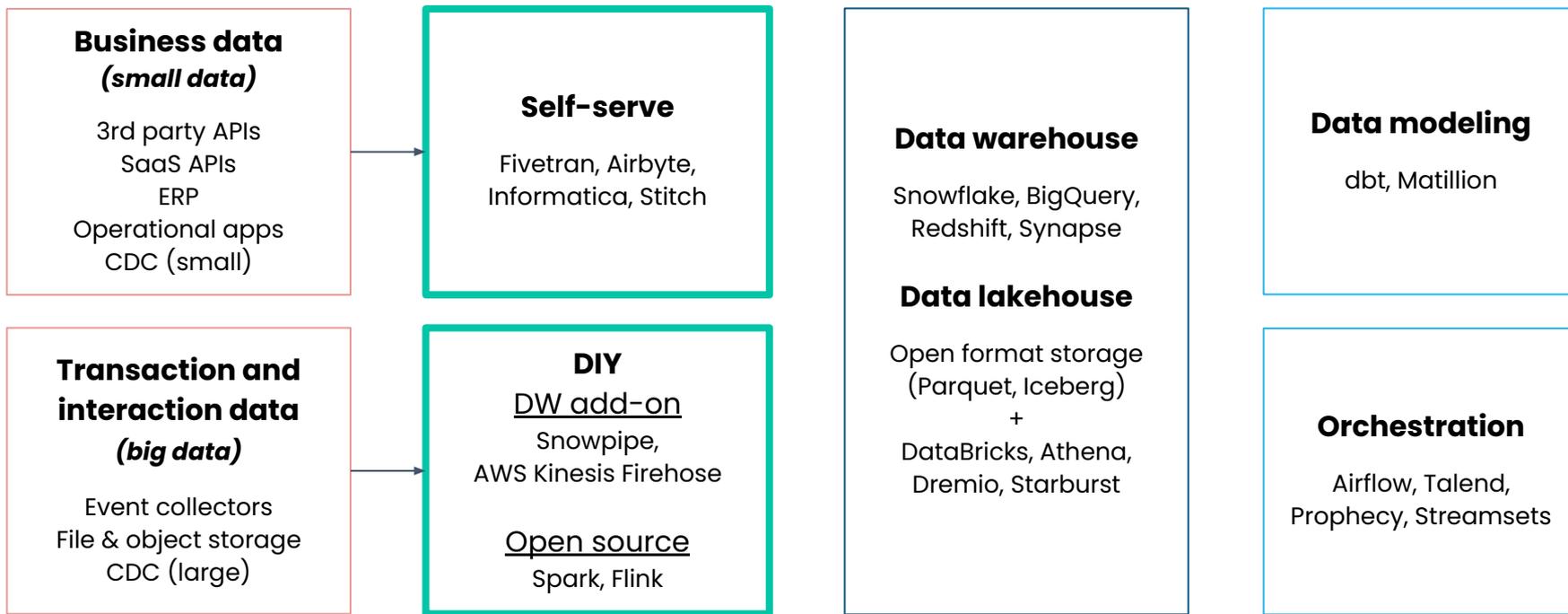
8% of customers	>1PB/mon
26% of customers	>100TB/mon
54% of customers	>10TB/mon
82% of customers	>1TB/mon

# The gap

1. Market map
2. Alternatives for high-scale ingestion
3. The tooling gap for high-scale ingestion

# Ingestion in the modern data stack

The modern data stack (MDS) makes it simple for companies to become data-driven



# For high-scale, existing tools fall short

## Ingestion

### Self-serve

Fivetran, Airbyte,  
Informatica, Stitch

### DIY

#### DW add-on

Snowpipe,  
AWS Kinesis Firehose

#### Open source

Spark, Flink

### Challenge

NOT built or priced for high-scale

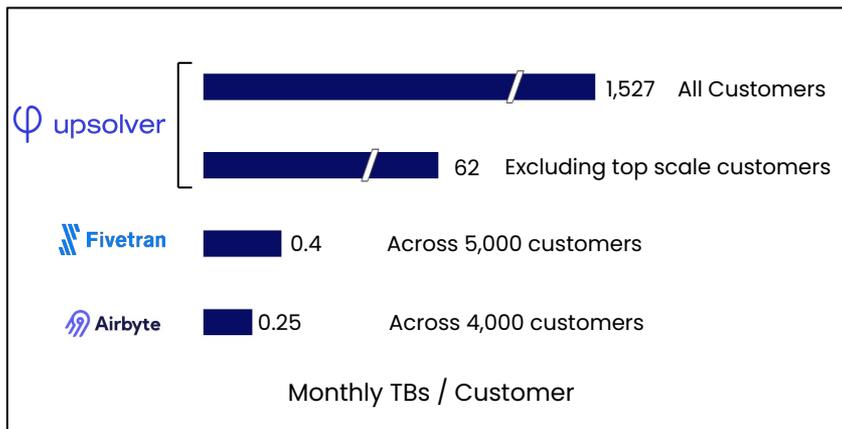
### Challenge

Data engineering bottleneck limits agility  
and causes many data quality incidents

# Self-serve ingestion tools are NOT built/priced for scale

Product focus is long tail connectors (SaaS, small CDC) with small business data

Today, customers use self-serve tools to process small scale data



## Fun fact

Upsolver's 2nd biggest customer has more data than all 9,000 Fivetran & Airbyte customers combined

At scale, self-serve tools' cost per month is 10X+ higher than Upsolver & DIY

Scale	Fivetran	Airbyte	Upsolver
1 TBs	\$22,000	\$5,000	\$4,225
10 TBs	\$50,000	\$20,000	\$6,250
100 TBs	\$500,000	\$200,000	\$26,500

## Assumptions:

- Append-only ingestion, 1KB per row
- Upsolver list price - \$4,000 + #TBs\*\$225

# DIY is a bottleneck to becoming data-driven

Counter to MDS ethos of pursuing higher value activities, not DIY data engineering



## Data initiatives orgs want to invest in

AI / ML

Personalization & segmentation

Data apps

Real-time insights

## Where time is spent with DIY

Integration between tools

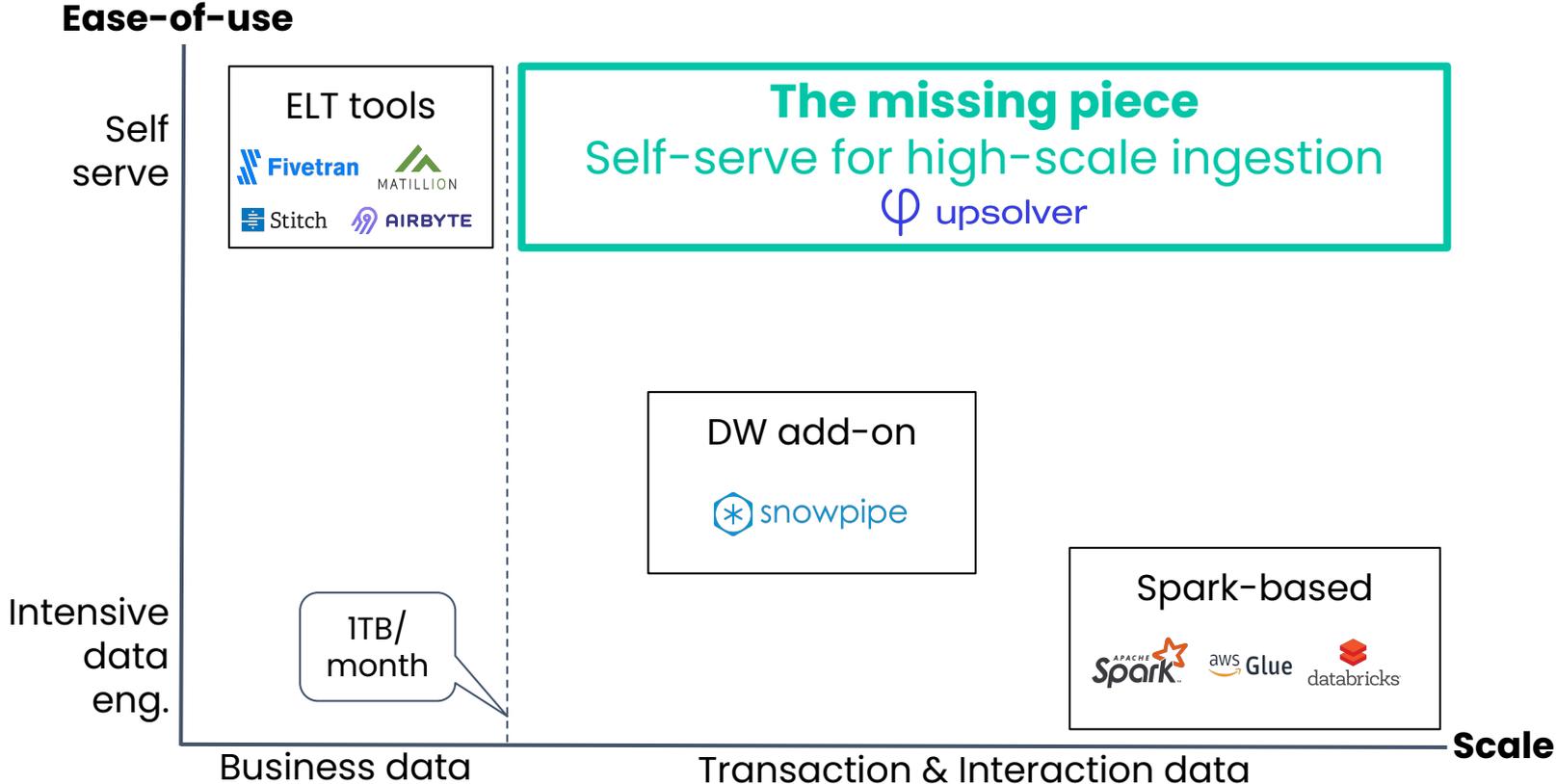
Scripting

Investigating quality problems

Fixing incorrect data

Evolving schema manually

# Upsolver fills the tooling gap for high-scale ingestion



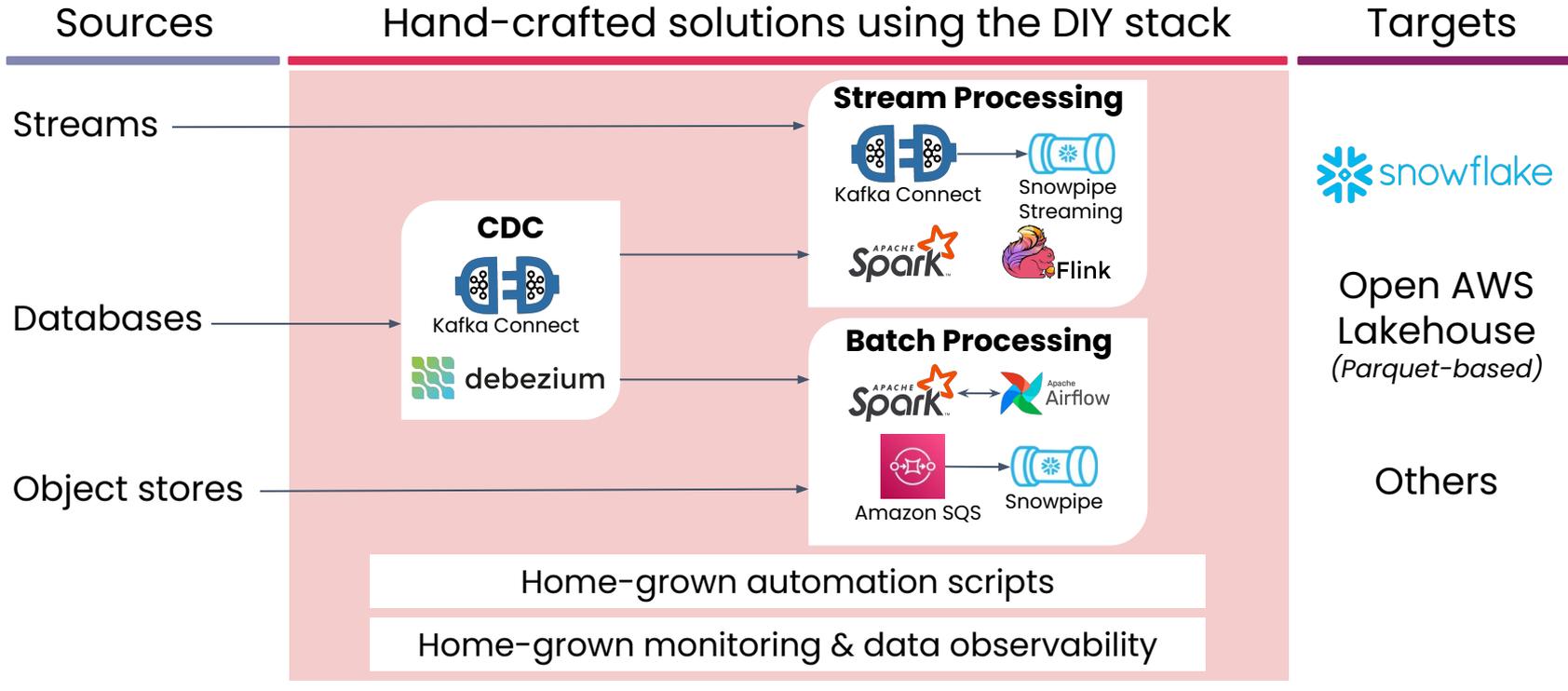
# Product

- The capabilities that solve the DIY complexity
  1. ETL using a single tool
  2. No-code & low-code dev experience
  3. Automatic prevention of most data disruptions
  4. Controlling quality – discover, resolve, prevent

# Upsolver uniquely bridges the DIY complexity gap

	DIY - Open source (Spark-based)	DIY - DW add-on (Snowpipe)	Self-serve (Upsolver)
Single tool for ETL	No	No	Yes
No-code and low-code dev experience	No	No	Yes
Automatic prevention of most data disruptions	No	No	Yes
Built-in tools for discovering, resolving and preventing data quality issues	No	No	Yes
<b>Outcome</b>	<b>Months to implement, Constant break-fix and ops headaches</b>		<b>Hours to implement, near-zero ops</b>

# (1) ETL using a single tool



 upsolver **1 tool instead of 3-5 moving pieces**

# (2) No-code & low-code dev experiences

**1** Set up your source

- Use an existing connection
  - upsolver\_kafka\_samples
  - upsolver\_kafka\_samples tested successfully
- Create a new connection

Select a topic to ingest: user\_info

Select the source events content type: Automatic

Sample Events

```
[{"id": "0", "data": {"$event_time": 1683531485840}}, {"id": "1", "data": {"$event_time": 1683531485840}}, {"id": "2", "data": {"$event_time": 1683531485840}}
```

**2** Set up your target

- Use an existing connection
  - snowflake\_target\_conn
  - snowflake\_target\_conn tested successfully
- Create a new connection

Select a schema: DEMO

New table name: USER\_INFO

**3** Configure the ingestion

Name your job: user\_info\_to\_USER\_INFO

How often do you want to update the target? (writing interval): 1 Minute

Which events to ingest? Start from now

Prevent duplicate events (deduplication)

Select fields for deduplication key: user\_id

Deduplication window: 1 Minute

Schema configuration

As soon as the schema changes, the newly added columns will be added to the target table.

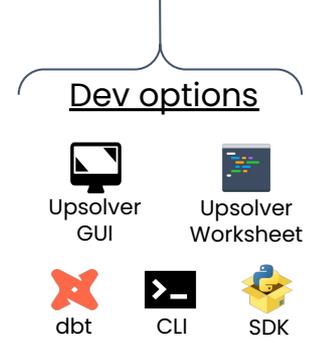
Search: address, credit\_card, first\_name, last\_name

**4** Review and run job

```
CREATE_SYNC_JOB upsolver_kafka_samples_to_upsolver_snowflake
CREATE_TABLE_IF_MISSING = true
START_FROM = BEGINNING
CONTENT_TYPE = AUTO
DEDUPLICATE_WITH = (COLUMNS = (orderid) WINDOW = 1 HOURS)
COLUMN_TRANSFORMATIONS = (customer_email = MD5(CAST(customer_email AS STRING)))
WRITE_INTERVAL = 1 HOURS
EVENT_TIME_COLUMN = UPSOLVER_EVENT_TIME
ADD_MISSING_COLUMNS = true
AS_COPY_FROM KAFKA upsolver_kafka_samples TOPIC = 'orders'
INTO SNOWFLAKE upsolver_snowflake.SUMMIT.ORDERS;
```

Previous Edit in Worksheet Run

Code is simple and therefore easy to understand, test and manage



load\_user\_events • Running Up to date 0 14 days ago

# (3) Automatic prevention of most data disruptions

## Death by a 1,000 silent data bugs

- Field type isn't valid in target
- Field name isn't valid in target
- Data type changed unexpectedly
- Numeric value casted incorrectly
- Field added/deprecated at the source
- Nested structs weren't flattened
- Events arrived late
- Unpredictable volume (outage / burst)

### Preventative

Automatically resolve conflicts & data bugs in real-time

Incident

\$

Ingestion

/\$

Staging

/\$

Transformation

/\$

Consumption

IMPACT

### DISASTER

Stale, incorrect, or corrupt data ingrained in analytics & ML products

### Extensive Fix

Multiple engineers  
Many pipelines  
Heavily used tables  
Multiple days

### Localized Fix

1 data engineer  
few pipelines  
few raw tables  
few hours



## Case studies

- VS self-serve ELT
- VS DIY add-on
- VS DIY open source

# Technology

1. Cloud-native architecture for high-scale processing
2. Key-value store for big time-series data
3. Intelligent landing zone for raw data
4. Easy-to-use CDC service, built on Apache Debezium

