

Streaming

Why should I care?

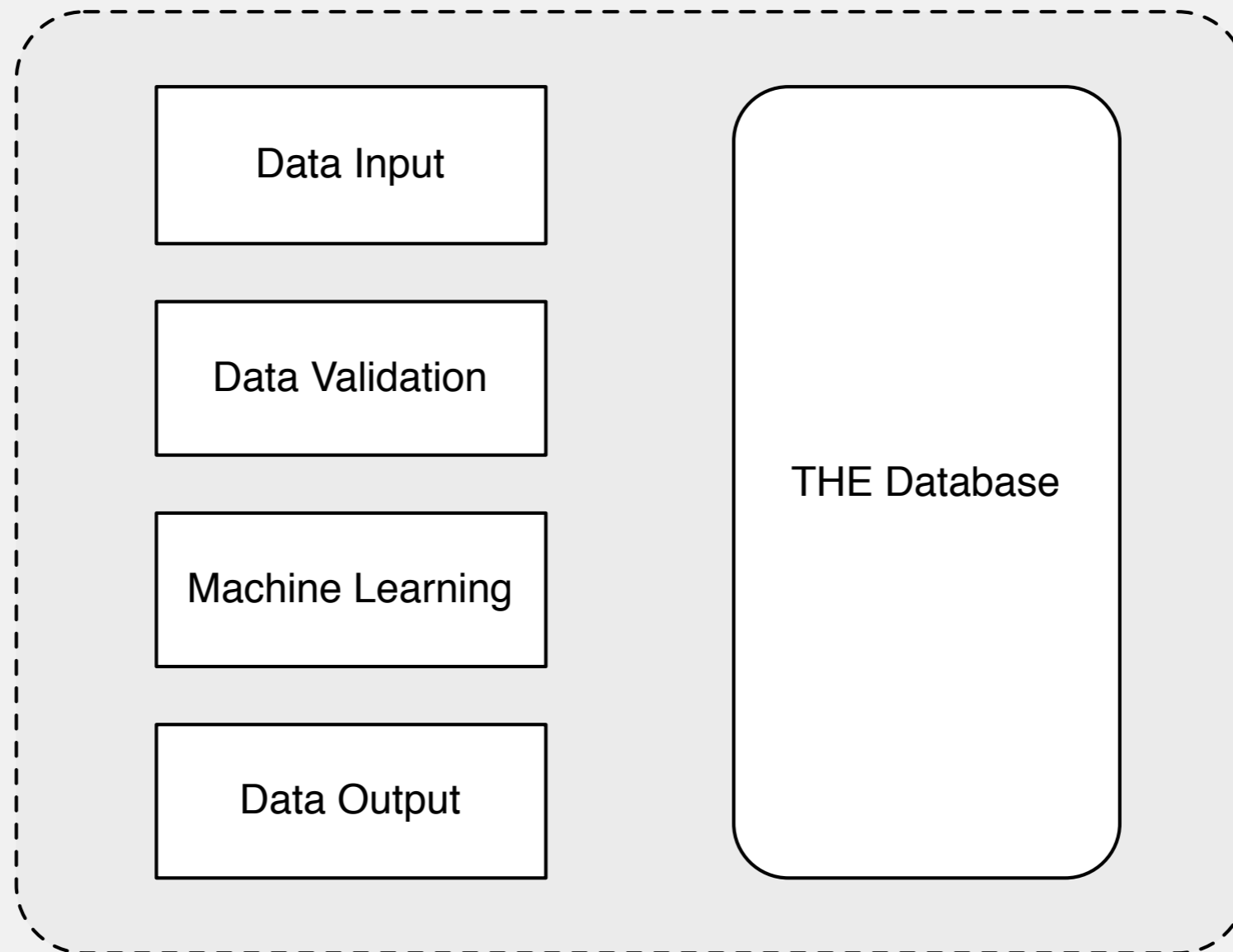
Christian Trebing
Blue Yonder GmbH
@ctrebing



Agenda

Motivation
Streaming Intro
Implementation
Challenges

Data Processing - The Monolith



Pain

Several teams are developing this application

- Customer desperately wants new feature in machine learning
- But the data validation team is in the midst of refactoring their database structure (‘will be finished in two weeks’)
- So you wait...

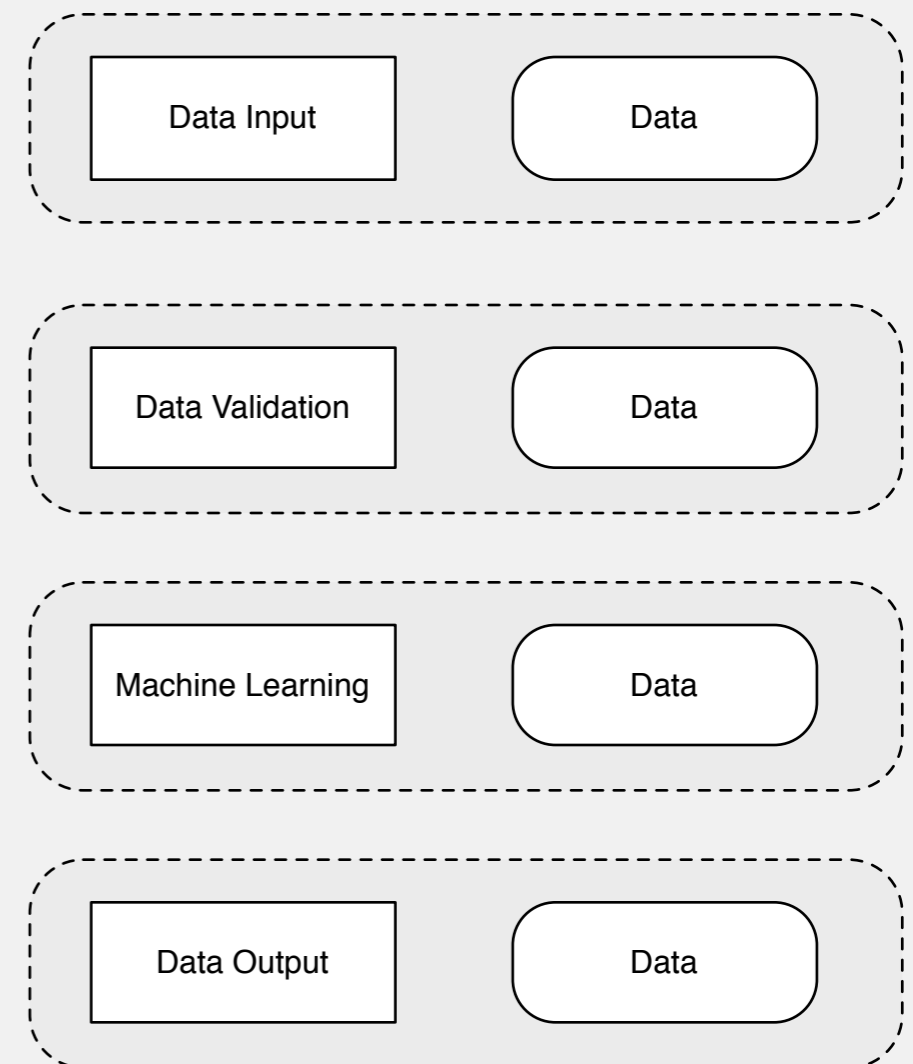
Could Microservices Help?

Great:

- No Dependency on single state
- Independent Development
- Independent Upgrades

Difficult:

- Too much data to transfer
- Too much data to store in each service



Are There Other Possibilities?

Streaming Intro

Databases and Streams

Same information in table and stream:

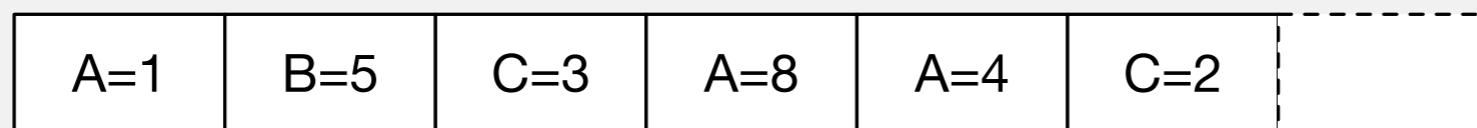
A=1	B=5	C=3	A=8	A=4	C=2	
-----	-----	-----	-----	-----	-----	--

A: 1	A: 1	A: 1	A: 8	A: 4	A: 4
	B: 5	B: 5	B: 5	B: 5	B: 5
		C: 3	C: 3	C: 3	C: 2

Why does it matter?

Different services can be in different states

- Each service can consume the stream in its own speed
- One service can be updated while the other runs



A: 1	A: 1	A: 1	A: 8	A: 4	A: 4
	B: 5	B: 5	B: 5	B: 5	B: 5
		C: 3	C: 3	C: 3	C: 2

Service 1
on index 3

Service 2
on index 5

Partitioned Streams

Sales stream, partitioned by location

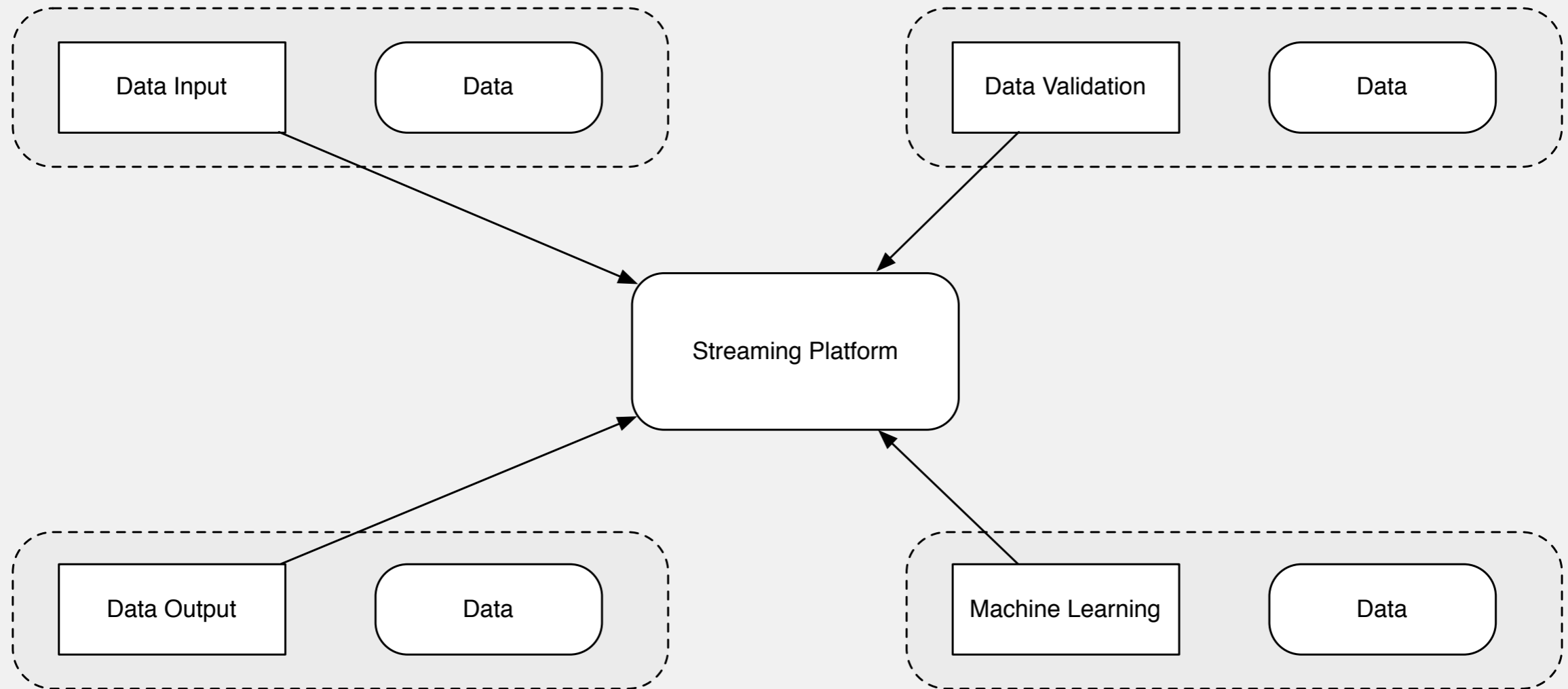
Each partition could be handled by diff. processors

location: Rimini product: Spaghetti sales_date: 2017-07-10 quantity: 5	location: Rimini product: Ravioli sales_date: 2017-07-10 quantity: 8	location: Rimini product: Pizza sales_date: 2017-07-11 quantity: 1	location: Rimini product: Spaghetti sales_date: 2017-07-11 quantity: 7
--	--	--	--

location: Bilbao product: Pizza sales_date: 2017-07-10 quantity: 3	location: Bilbao product: Ravioli sales_date: 2017-07-11 quantity: 9	location: Bilbao product: Spaghetti sales_date: 2017-07-11 quantity: 5
--	--	--

location: Karlsruhe product: Pizza sales_date: 2017-07-10 quantity: 8	location: Karlsruhe product: Ravioli sales_date: 2017-07-10 quantity: 3	location: Karlsruhe product: Ravioli sales_date: 2017-07-11 quantity: 7	location: Karlsruhe product: Spaghetti sales_date: 2017-07-11 quantity: 7
---	---	---	---

Data Processing - With Streaming



What did we gain?

Independent Development

Independent Upgrade

Scalability

Did we throw out databases completely?

- Let's see...

Is it magic?

No, it's a tradeoff:

- A database is so powerful: ACID guarantees, SQL language. You can do nearly everything
 - This comes at a price
 - Dependency on single state
 - Scaling is hard

So let's more think what we really need

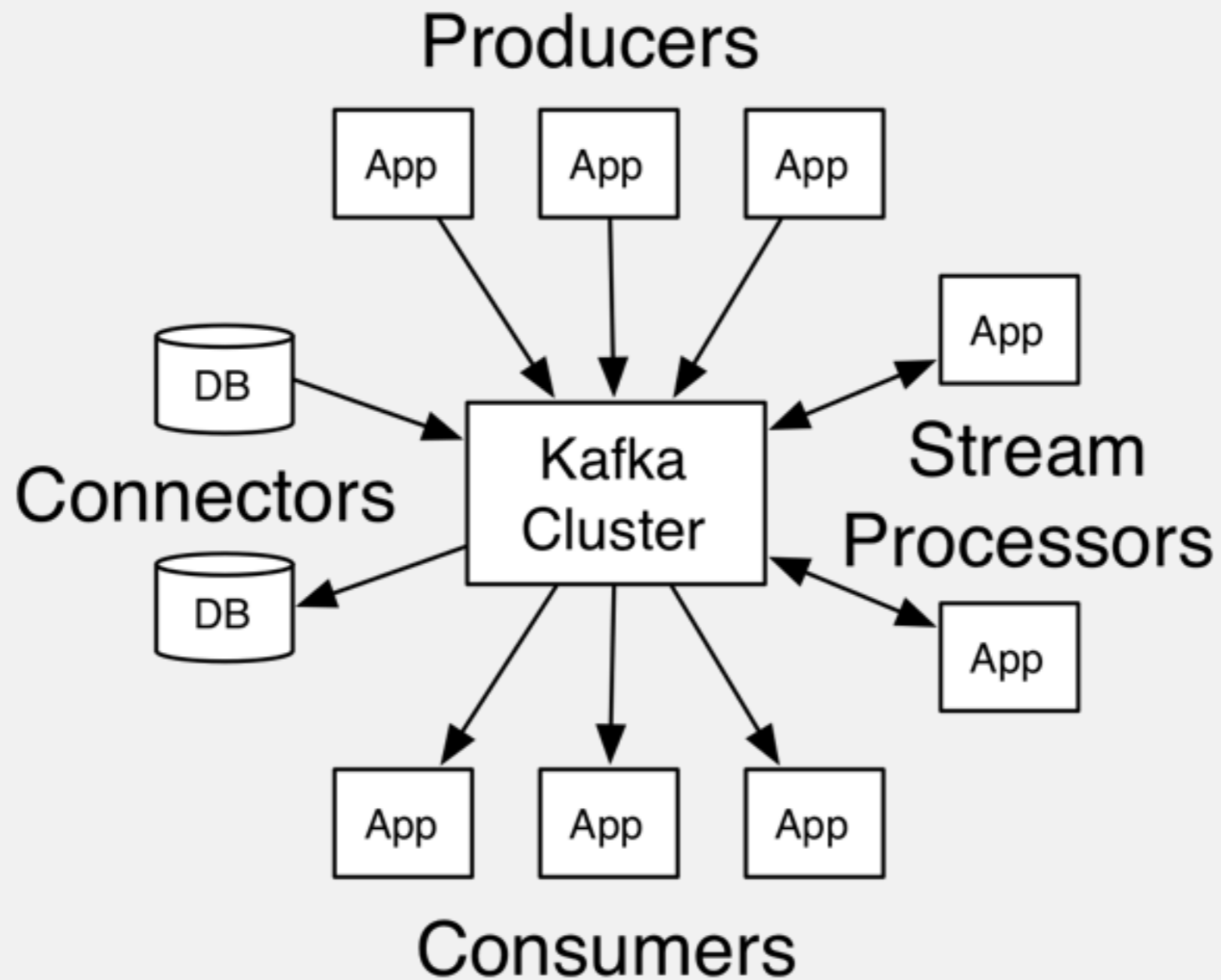
What do we loose?

Database	Stream
ACID	Ordering on stream partition
SQL Queries	Service is responsible of keeping its state

You have to decide whether you can live with that.

Implementation

Apache Kafka



Kafka Clients in Python

- pykafka, python-kafka, confluent-kafka-client
- Nice comparison has been done here:
 - <http://activisiongamescience.github.io/2016/06/15/Kafka-Client-Benchmarking/>
- Most performant currently is confluent-kafka-client
 - Uses the c library librdkafka

Producer

```
1 import json
2 from confluent_kafka import Producer
3
4 p = Producer({'bootstrap.servers': 'localhost:9092'})
5 data = {
6     'location': 'Rimini', 'product': 'Ravioli',
7     'sales_date': '2017-07-10', 'quantity': 5
8 }
9 p.produce('sales_input', json.dumps(data))
10 p.flush()
```

Consumer

```
1 from confluent_kafka import Consumer, KafkaError
2 import uuid
3
4 c = Consumer({'bootstrap.servers': 'localhost:9092', 'group.id': uuid.uuid1(),
5             'default.topic.config': {'auto.offset.reset': 'smallest'}})
6 c.subscribe(['sales_input'])
7 running = True
8 while running:
9     msg = c.poll()
10    if not msg.error():
11        print('Received message: %s' % msg.value().decode('utf-8'))
12    elif msg.error().code() != KafkaError._PARTITION_EOF:
13        print(msg.error())
14        running = False
15 c.close()
```

Apache Avro

Data Serialization, Enabling Schema Evolution



Clearly defined schema with:

- Schema evolution
- Schema registry

Writer's schema

Data Type	Field Name
string	location
string	product
string	sales_date
int	quantity

Reader's schema

Data Type	Field Name
string	location
string	product
string	sales_date
int	quantity
int, default=0	delivery_id

Avro Schema

```
1 import json
2 from confluent_kafka import avro
3
4 value_schema = avro.loads(json.dumps({
5     "name": "sales_input_value",
6     "type": "record",
7     "fields": [
8         {"name": "location", "type": ["string"]},
9         {"name": "product", "type": ["string"]},
10        {"name": "sales_date", "type": ["string"]},
11        {"name": "quantity", "type": ["int"]}
12    ]
13 })))
```

AvroProducer

```
1  from confluent_kafka import avro
2  from confluent_kafka.avro import AvroProducer
3  from schema import key_schema, value_schema
4
5  avroProducer = AvroProducer({
6      'bootstrap.servers': 'localhost:9092',
7      'schema.registry.url': 'http://localhost:8081'},
8      default_key_schema=key_schema, default_value_schema=value_schema)
9
10 key = {'location': 'Rimini'}
11 value = {
12     'location': 'Rimini',
13     'product': 'Ravioli',
14     'sales_date': '2017-07-10',
15     'quantity': 5
16 }
17 avroProducer.produce(topic='sales_input', value=value, key=key)
18 avroProducer.flush()
```

AvroConsumer

```
1  import uuid
2  from confluent_kafka import KafkaError
3  from confluent_kafka.avro import AvroConsumer
4  from confluent_kafka.avro.serializer import SerializerError
5
6  c = AvroConsumer({
7      'bootstrap.servers': 'localhost:9092', 'group.id': uuid.uuid1(),
8      'schema.registry.url': 'http://localhost:8081',
9      'default.topic.config': {'auto.offset.reset': 'smallest'}})
10 c.subscribe(['sales_input'])
11 running = True
12 while running:
13     try:
14         msg = c.poll(10)
15         if msg:
16             if not msg.error():
17                 print(msg.value())
18             elif msg.error().code() != KafkaError._PARTITION_EOF:
19                 print(msg.error())
20                 running = False
21     except SerializerError as e:
22         print("Message deserialization failed for %s: %s" % (msg, e))
23         running = False
24
25 c.close()
```

Example: Data Validation

Separate valid and invalid sales records

```
1 consumer.subscribe(['sales_input'])
2
3 while True:
4     msg = self.consumer.poll(1)
5     if msg:
6         key = msg.key()
7         value = msg.value()
8         if sales_record_valid(value):
9             producer.produce(topic='sales_validated', value=value, key=key)
10        else:
11            producer.produce(topic='sales_error', value=value, key=key)
```


Additional Processors

Need to evolve your application:

- Add processors for evaluation topics
- Try new variant of validation logic

database

- remember processing state, same for each processor

streaming

- offset in each processor, can work independently

Challenge

- Machine Learning Still Batch

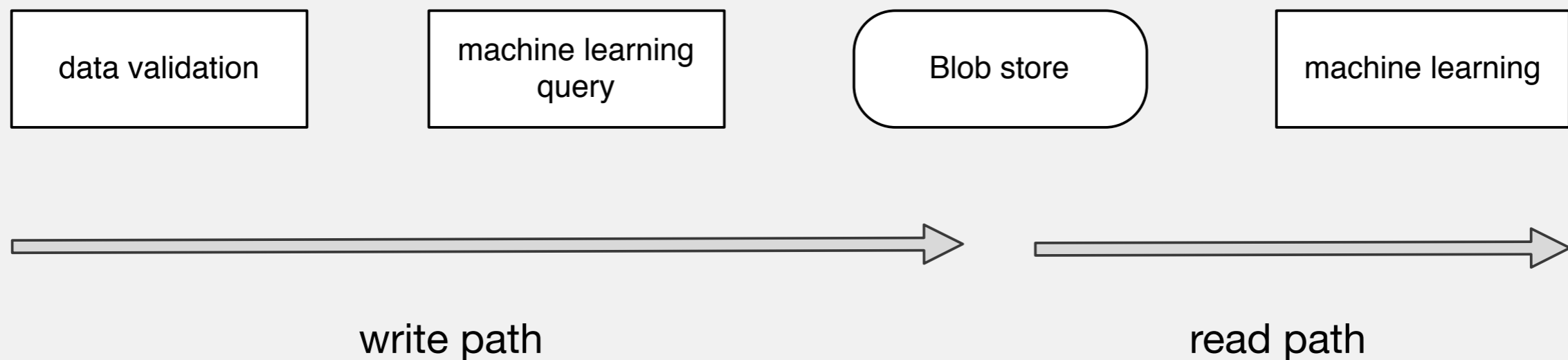
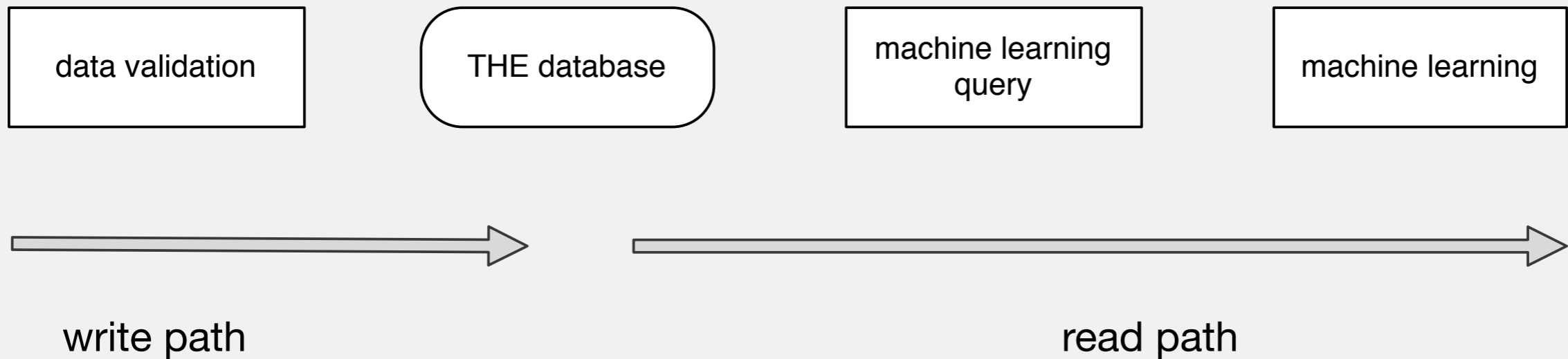
How to get the input data?

Remember: There is no possibility to query a stream. Somewhere all this data needs to be. Options:

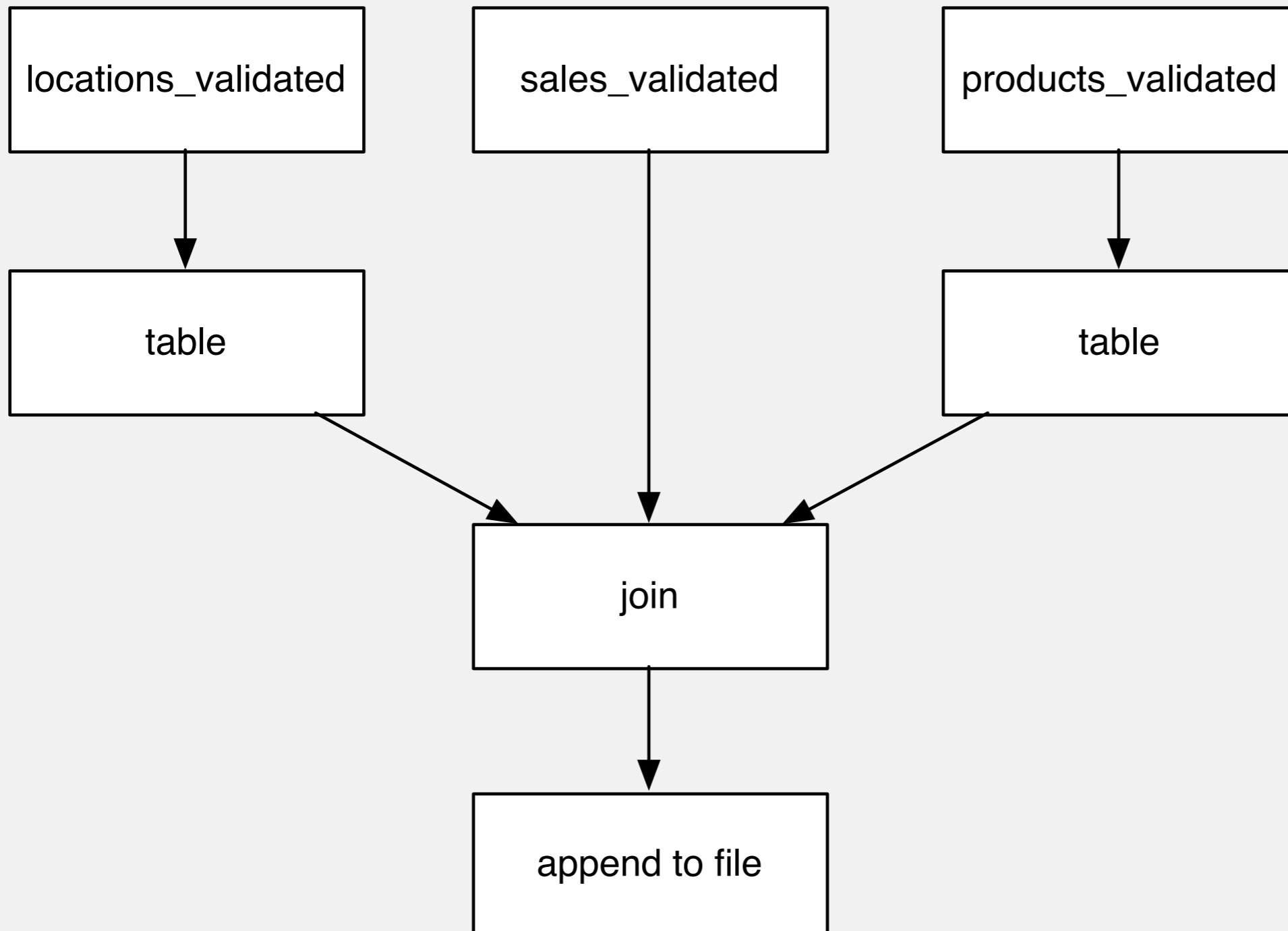
- Memory of the service
- Serving database
- Blob store

Yes, that's duplication. We'll have to live with it.

Write Path / Read Path



Machine Learning Input Data



Challenge

- State in Processors

State - Nightmare of every distributed systems engineer

Streaming: Data just rushes through

Why do we need state?

- Time window processing
- Data you want to join with

Formerly, the database did it for you

State - Some Challenges?

Failure of a processor

Scaling

State in Stream Processors

- Possible Solutions

- Just keep in memory. Reprocess stream to warm up
- Each processor to keep its own db
- Save condensed in stream
- Get it from other service

Frameworks exist in other languages:

- for example: Kafka Streams, Apache Samza

Up to yesterday: none in python. Then heard about

<https://github.com/wintoncode/winton-kafka-streams>

Summary

- You have more options for your data processing applications than you might have thought
- As always, there are some tradeoffs
- You know the challenges

Questions?

Blue Yonder

Best decisions, delivered daily

Blue Yonder GmbH
Ohiostraße 8
76149 Karlsruhe
Germany
+49 721 383117 0

Blue Yonder Software Limited
19 Eastbourne Terrace
London, W2 6LG
United Kingdom
+44 20 3626 0360

Blue Yonder Analytics, Inc.
5048 Tennyson Parkway
Suite 250
Plano, Texas 75024
USA