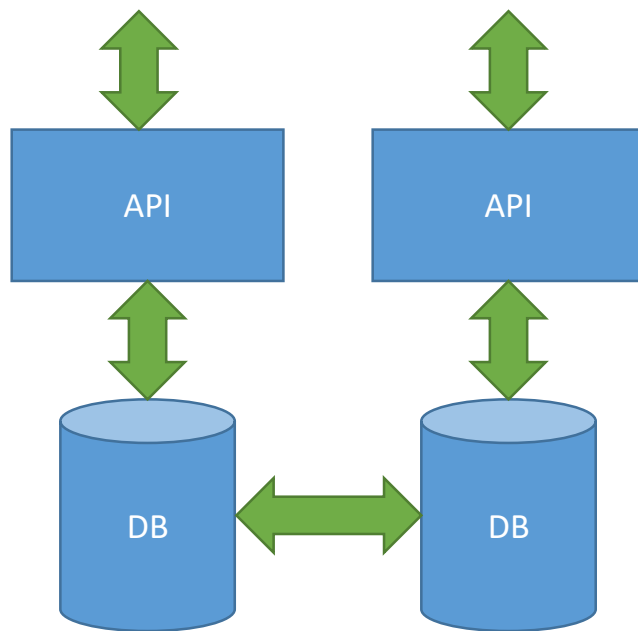


Polyglot Persistence in Azure

Ken Seier

Why Polyglot Persistence

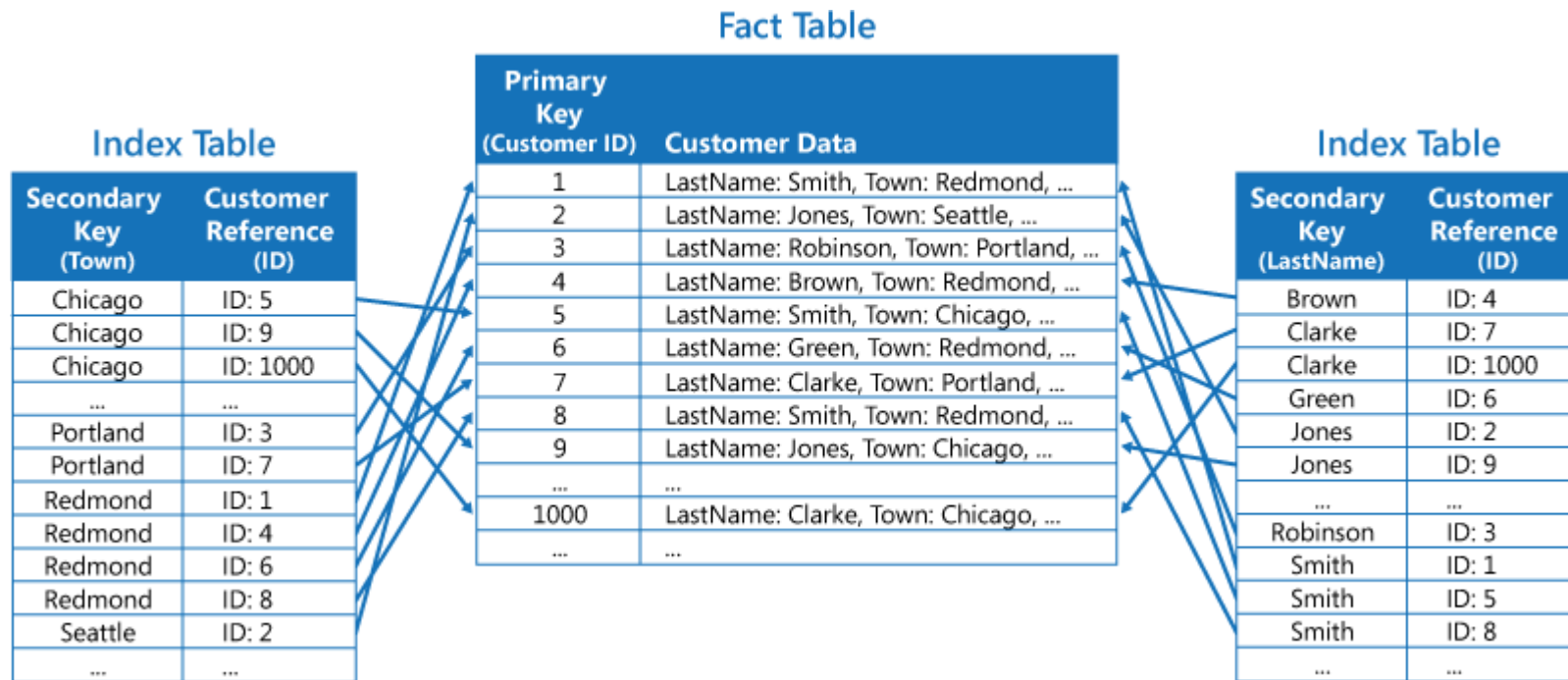
- Polyglot Persistence is storing data in multiple stores, often with different formats
- Term coined in 2006
- So much more than microservices
- It's about maintaining performance for low total cost of ownership (TCO) *in large systems*



Performance hockey stick



Familiar Polyglot pattern



Polyglot challenges



- Data consistency
- Transactional integrity
- Implementation complexity
- Creates data silos

A photograph of a business meeting with a pink overlay. Three people are gathered around a table, looking at a laptop and some papers. One person is pointing at the laptop screen, another is holding a pen over a document, and a third is holding a cup of coffee. The word "Workloads" is written in white text in the center of the image.

Workloads

Transactional workloads

- Have clearly defined use patterns
- Focused on create, read, update, delete (CRUD)
- Affect a very limited volume of data, ideally single record
- Ideally write-few/read-few but often write-few/read-many



Slicing workloads



- Have clearly defined use patterns
- Covers large numbers of records
- Slices (sub-selects) and aggregates data
- Commonly used for reporting and metrics
- Can be significantly optimized with indexing and caching

Exploratory workloads

- Has no defined pattern
- Often cover very large amounts of data
- Thwarts optimization



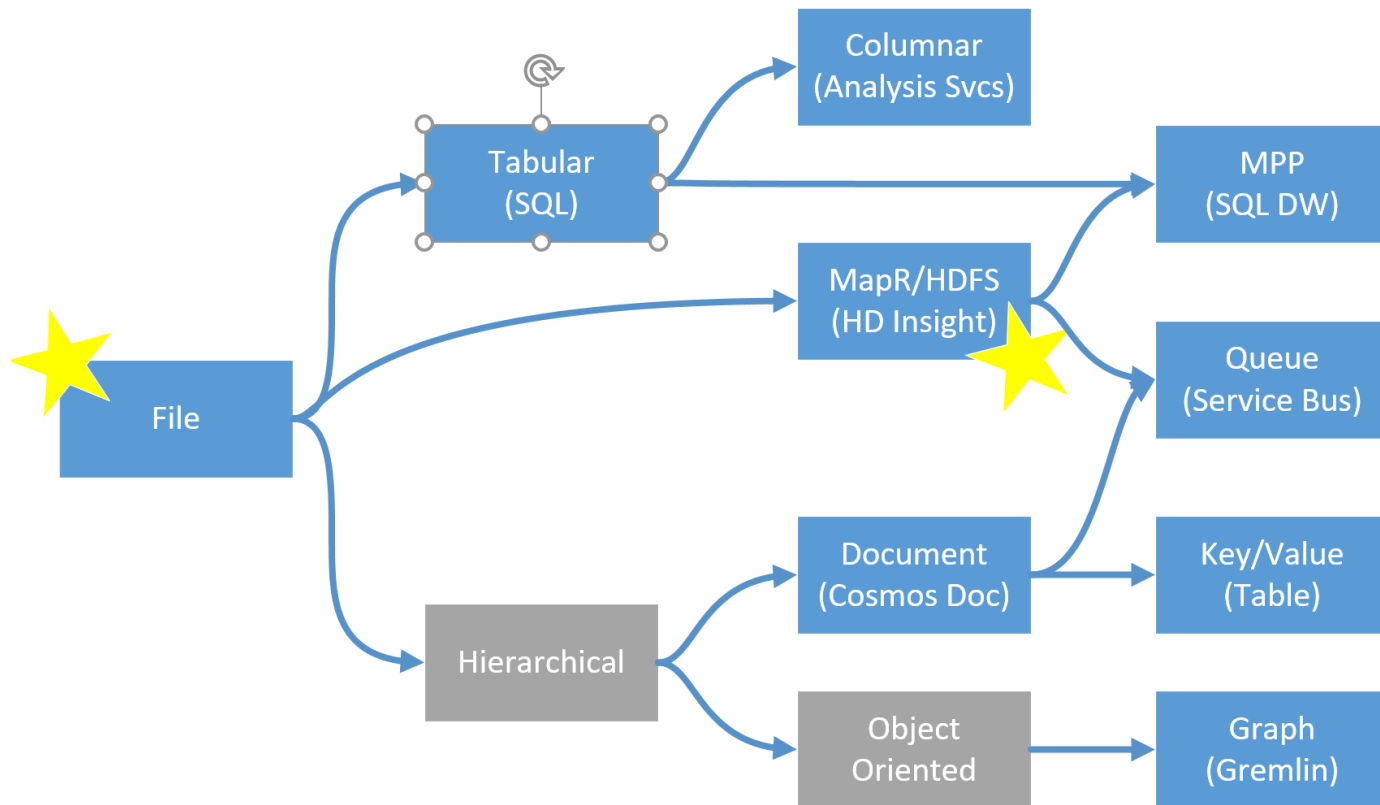
Different types of workloads in a single store



A photograph of a business meeting with a magenta color overlay. Three people are gathered around a table, looking at a document with charts. A laptop and a cup of coffee are also on the table. The text 'Data Stores' is centered in white.

Data Stores

Data store family tree



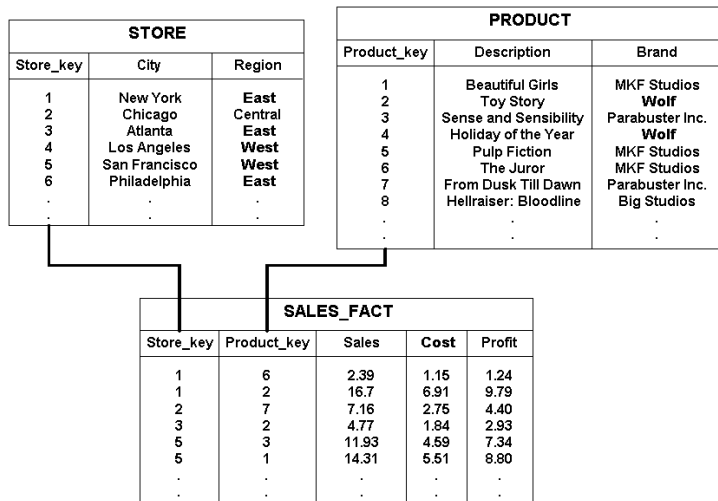
Many choices

- Driven by:
 - Backwards compatibility
 - Available skill
 - Future strategy

PAIN



Tabular stores



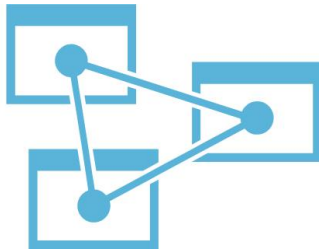
- Stores data in rectangular tables of rows with a fixed set of columns
- Tables are related via keys
- Enables strong normalization (deduplication) of data
- Good for transactional and slicing workloads
- OK for exploratory workloads



Columnar stores

- Stores data in tables of columns with each column fully indexed
- Relies heavily on memory and caching for performance
- Excellent at slicing workloads
- OK at transactional workloads
- Poor at exploratory workloads

Record #	Name	Address	City	State
0003623	ABC	125 N Way	Cityville	PA
0003626	Newburg	1300 Forest Dr	Troy	VT
0003647	Flotsam	Industrial Pkwy	Springfield	MT
0003705	Jolly	529 S 5th St.	Anywhere	NY



Document stores

```
{
  "Name": "Yogurt Depot",
  "Id": 1,
  "Revenue": 2000,
  "Cost": 100,
  "Category": [ "dessert", "food", "yogurt" ],
  "Visits": [
    {
      "day": "Mon",
      "visit_count": 300
    },
    {
      "day": "Tue",
      "visit_count": 700
    }
  ]
},
```



- Document defines a single record
- Document can contain just about any configuration of text-based data (JSON)
- Excellent at single-document transactional workloads
- Poor at slicing and exploratory workloads

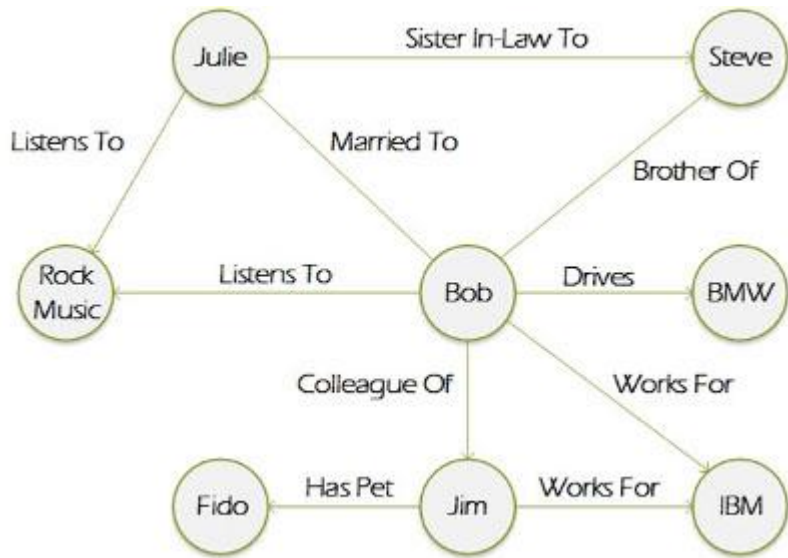
Key/value stores

- Stores binary blob value referenced by a unique key
- Blobs can be anything, but are often ragged column arrays
- Data often partitioned by key
- A smart key can improve performance by acting as an index
- Excellent at single record transactional workloads
- Poor at slicing and exploratory workloads

Key	Value
K1	AAA,BBB,CCC
K2	AAA,BBB
K3	AAA,DDD
K4	AAA,2,01/01/2015
K5	3,ZZZ,5623



Graph stores

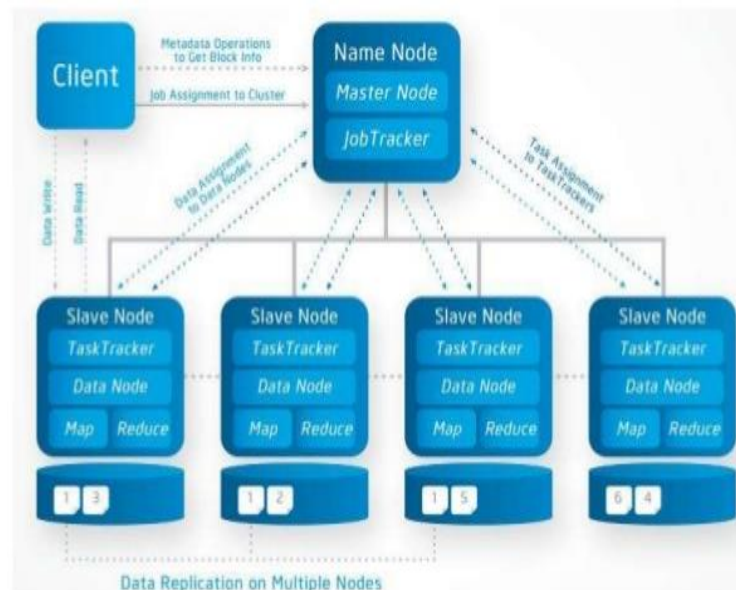


- Stores data as nodes and node relationship (edges)
- Excellent at network-type workloads
- Poor at all other workloads

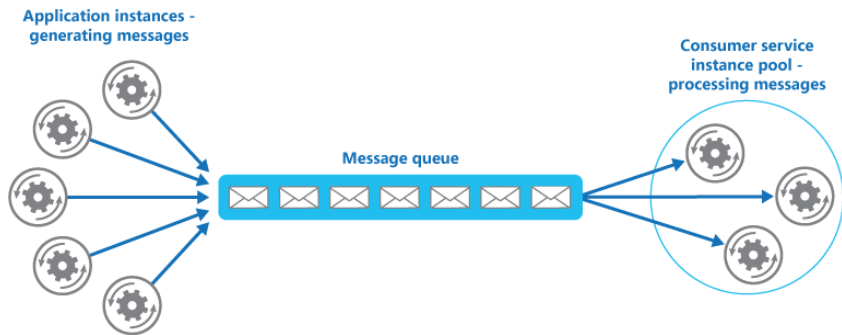


MapR/HDFS stores

- Based on a “head” node controlling multiple “worker” nodes
- Can accept any type of binary data
- Extremely flexible
- Excellent at exploratory workloads
- Good at transactional and slicing workloads over large amounts of data



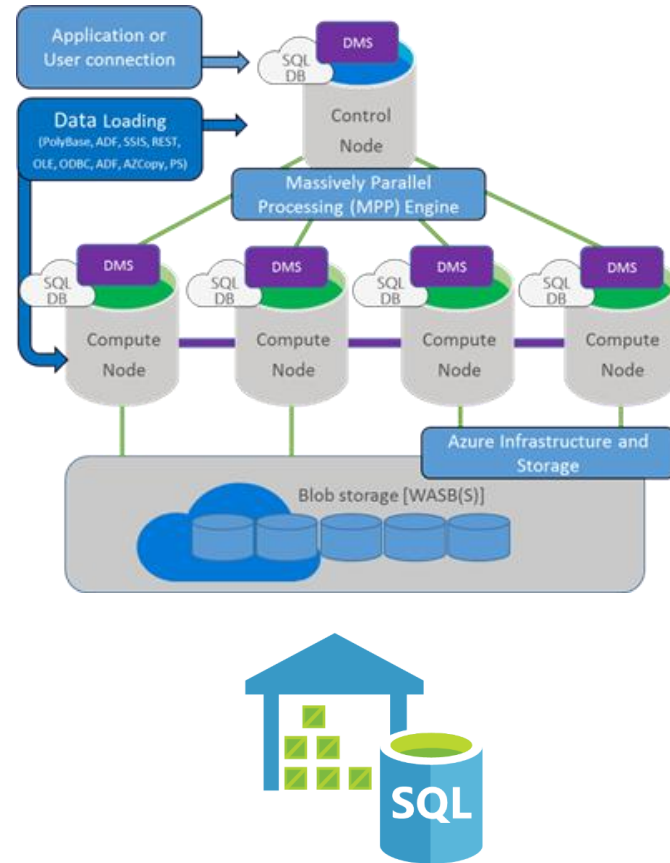
Queue stores



- Stores messages (documents) in ordered format
- Optimized for ordered consumption messages
- Multiple patterns for handling messages: read-only, consume, reservation, etc.
- Excellent at some transactional workloads
- Terrible at slicing and exploratory workloads

MPP stores

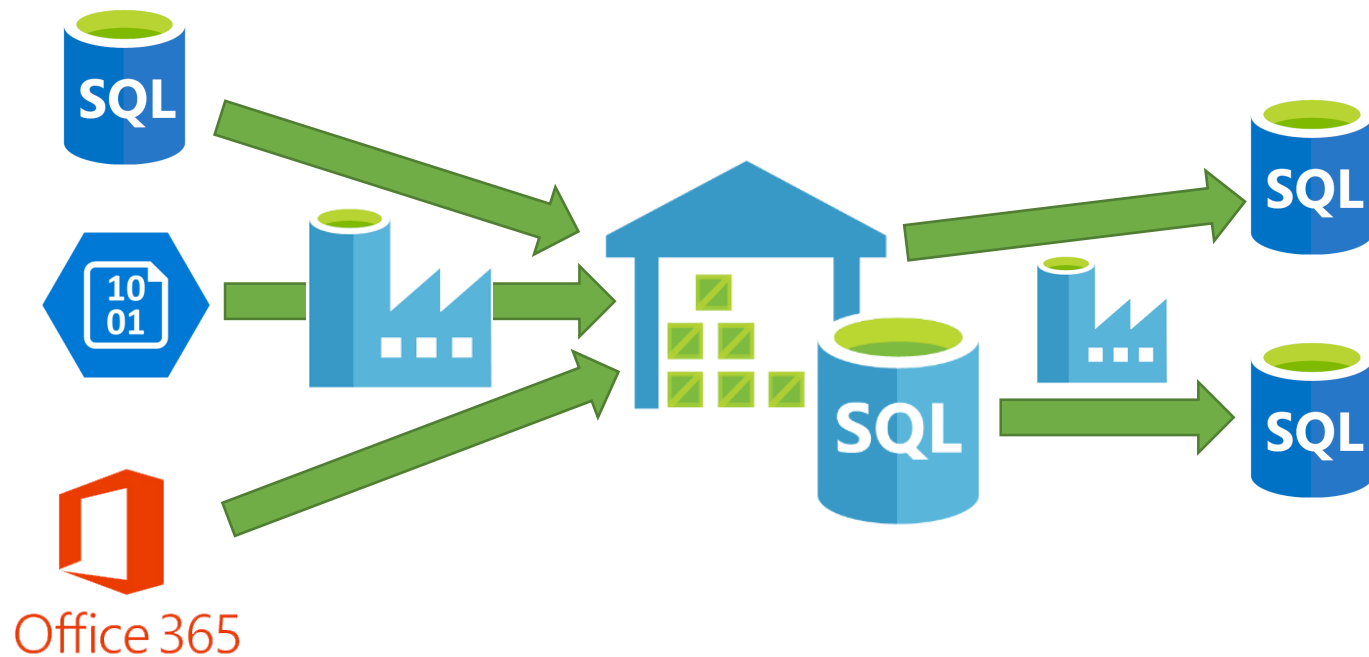
- Stores tabular data across multiple horizontal partitions (shards)
- Distributes large tables across shards using partition key
- Replicates small tables across shards
- Excellent at slicing workloads over large amounts of data
- OK at transactional workloads
- Poor at exploratory workloads



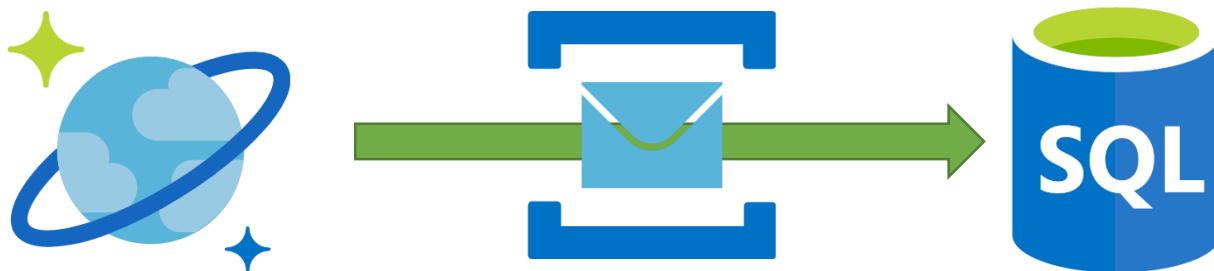
A photograph of a group of people in a meeting, overlaid with a semi-transparent purple filter. The scene shows several individuals gathered around a table, looking at documents and a laptop. One person is pointing at a document, while another is writing. A cup of coffee is visible in the foreground. The text 'Polyglot Patterns' is centered in white.

Polyglot Patterns

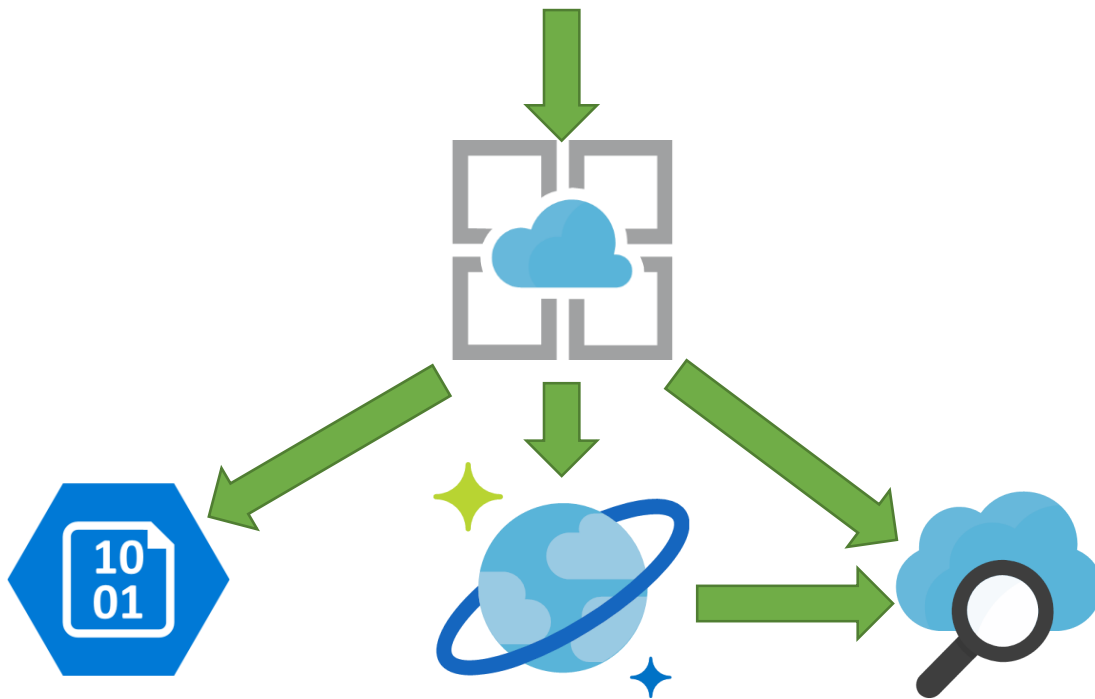
OLTP/OLAP



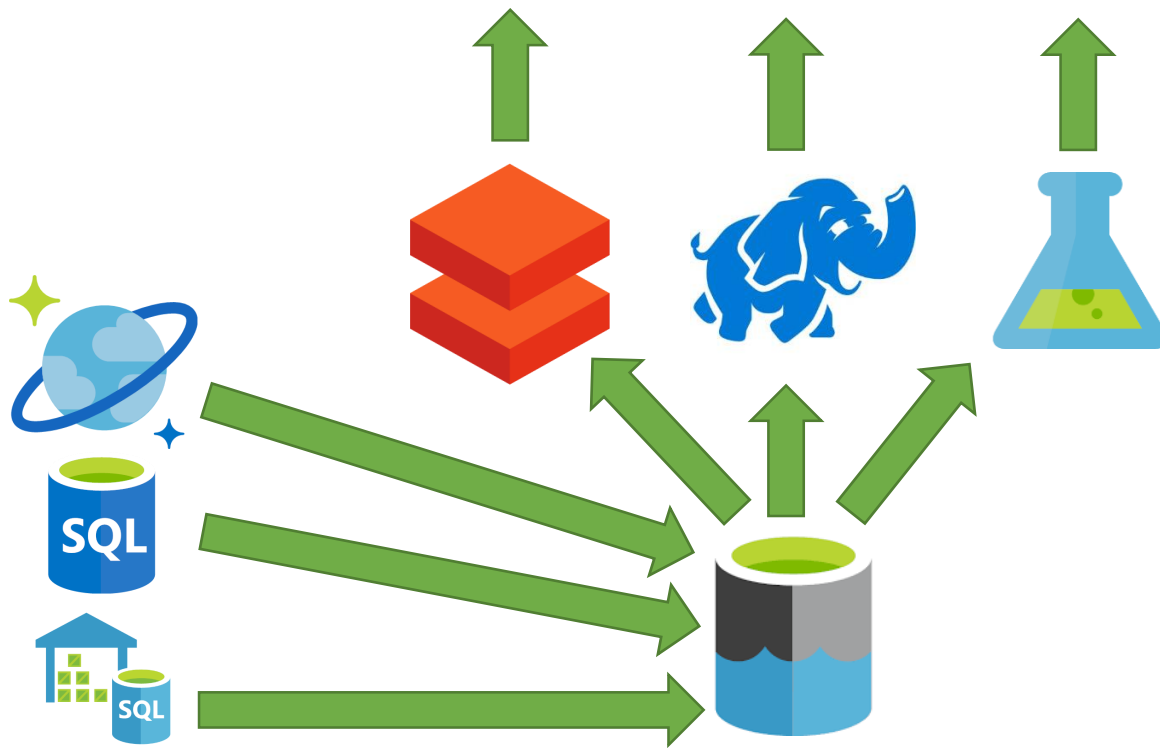
Dependent systems



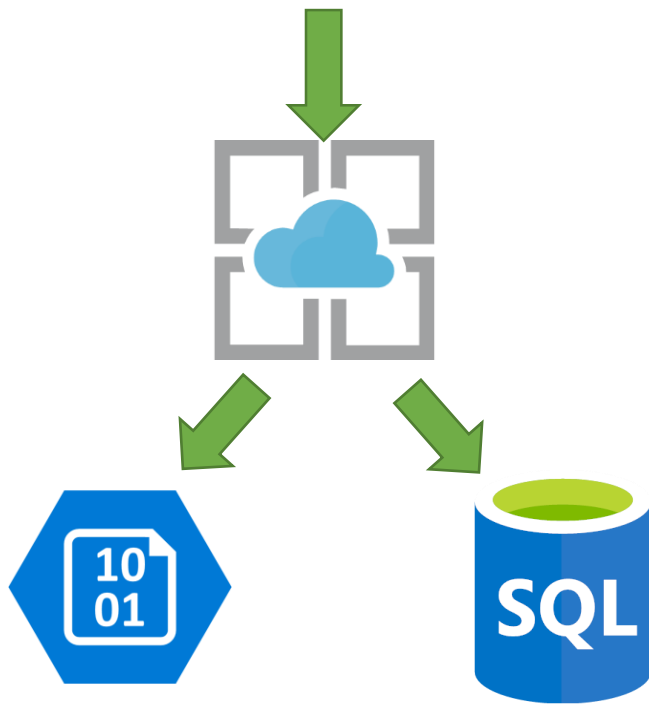
Indexing systems



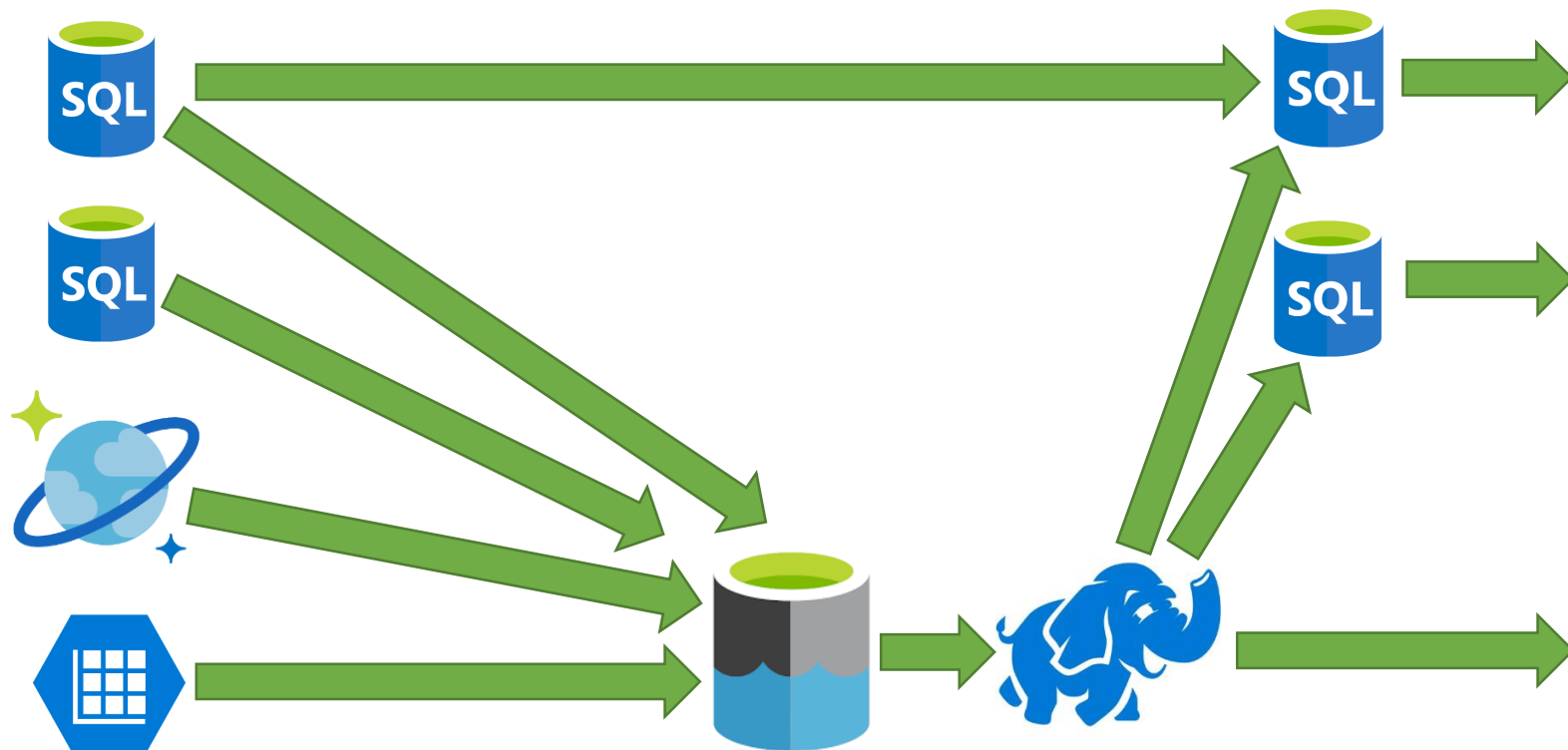
Exploratory systems



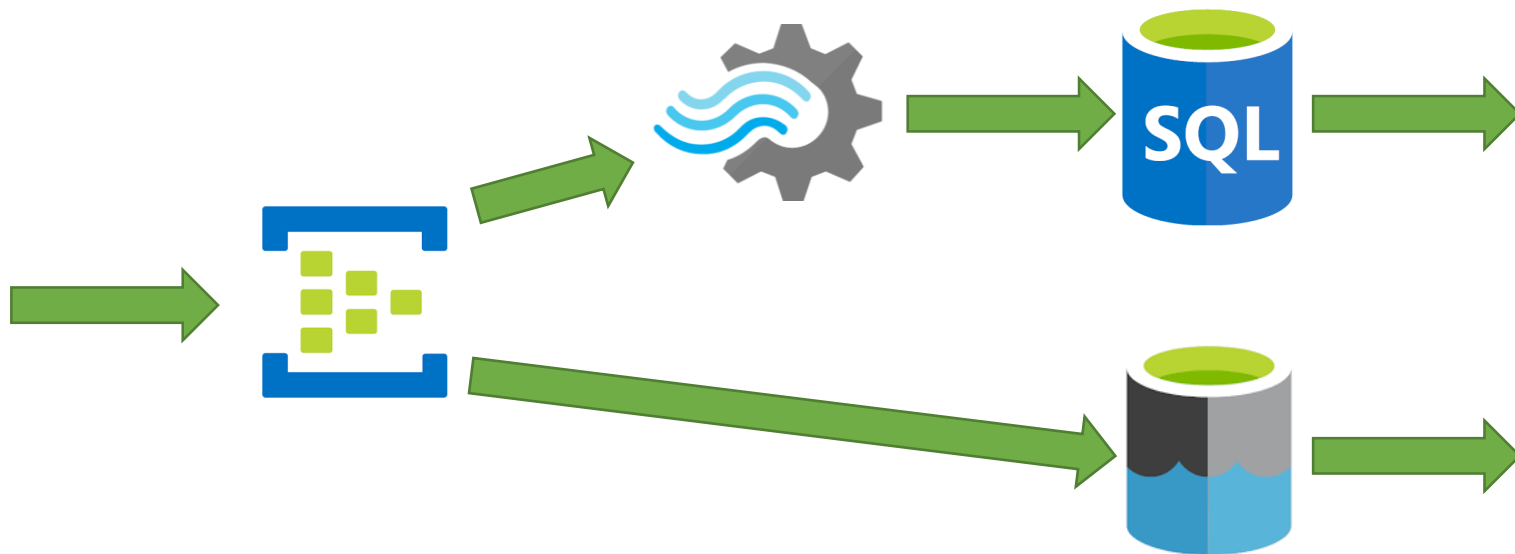
Disparate data types



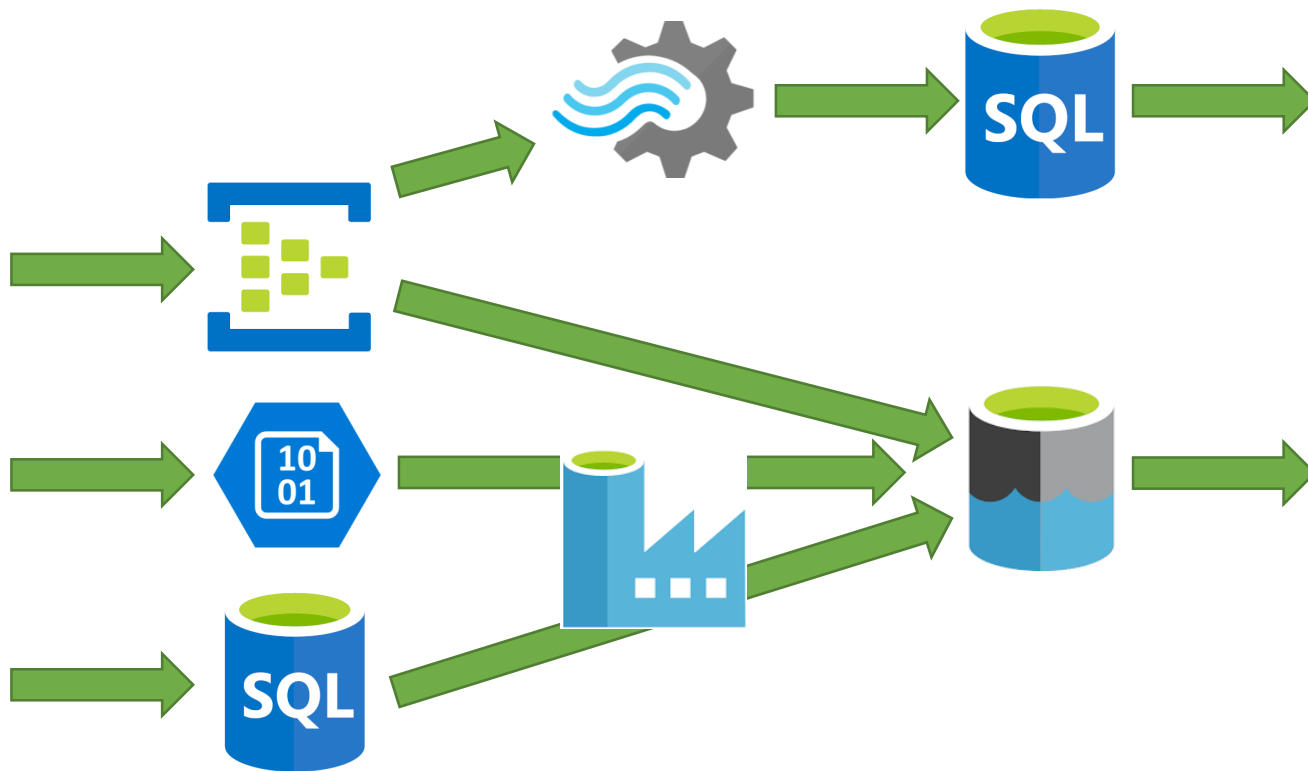
Modern data warehouse



Streaming systems (Kappa architecture)



Lambda architecture



A photograph of a business meeting with a pink overlay. Three people are gathered around a table, looking at a document. A laptop is open in the background, and a cup of coffee is in the foreground. The word "Summary" is written in white text in the center of the image.

Summary

Bringing it together

- Polyglot persistence is aligning data stores with workloads to improve TCO. This can mean:
 - Performance costs
 - Platform/hardware/software costs
 - Development velocity
- Trades implementation complexity for TCO benefits
- Only needed in higher-complexity systems (large code base, large data, high performance needs, etc.)
- Lots of good choices, even more bad choices
- SQL isn't going anywhere



Thank You

A BIG thank you to the 2018 Global Sponsors!

