

Lessons Learned Running Big Data on AWS

Michael Rhee, Architect – Gogo Data Warehouse and Analytics

August 1, 2018





Gogo is the inflight internet company.

AVIATION-CENTRIC

RELENTLESSLY INNOVATIVE

PERFORMANCE OBSESSED

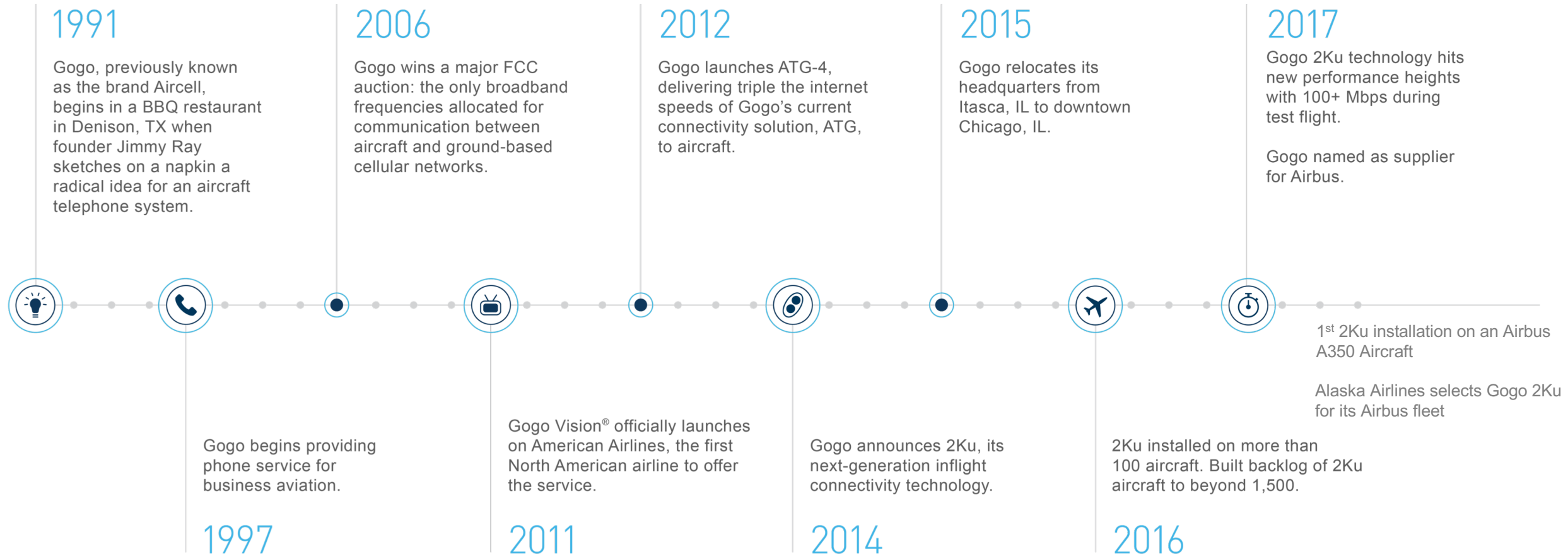
About Gogo



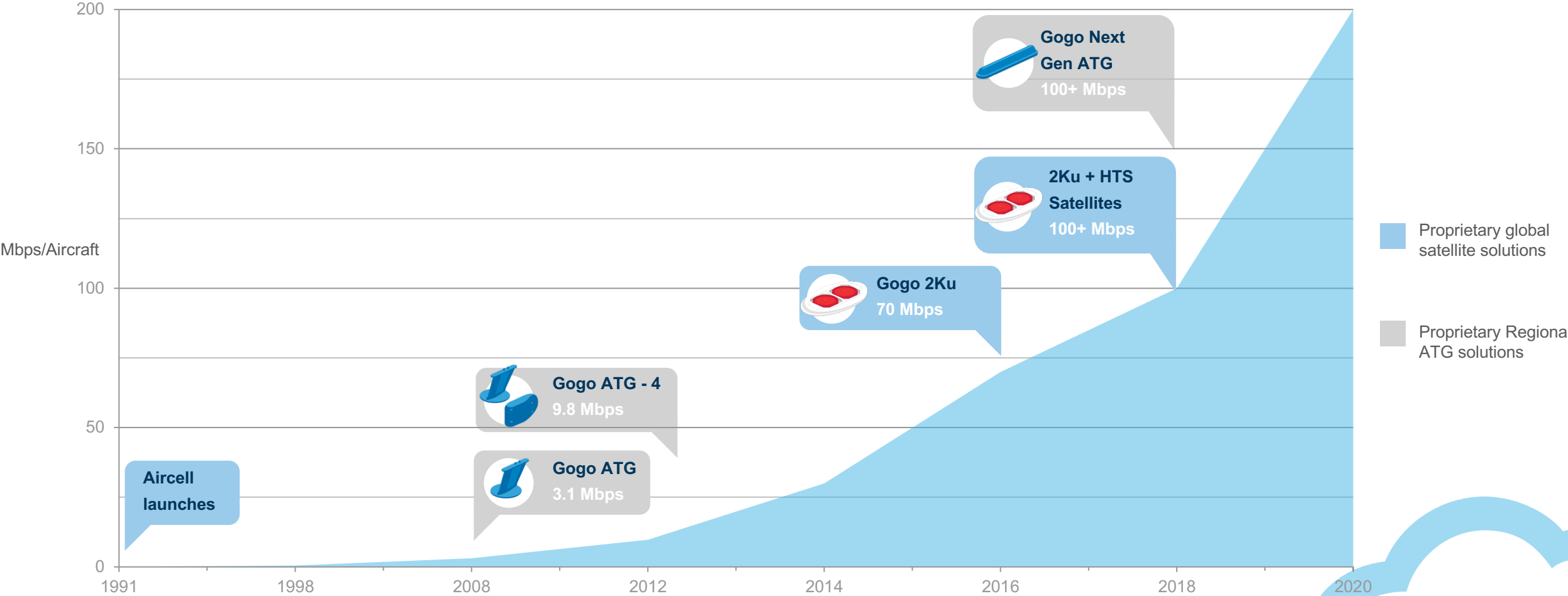
- 7,200+ Connected Aircraft
- 21 commercial airline partners
- Partnerships with the largest fractional ownership operators in business aviation
- 2,000+ aircraft awards for 2Ku, Gogo's latest commercial aviation technology
 - 19 airlines committed to 2Ku
- 120M+ connectivity sessions to date
- 140,000+ sessions/day
- 11,000+ flights/day



Gogo History at a Glance

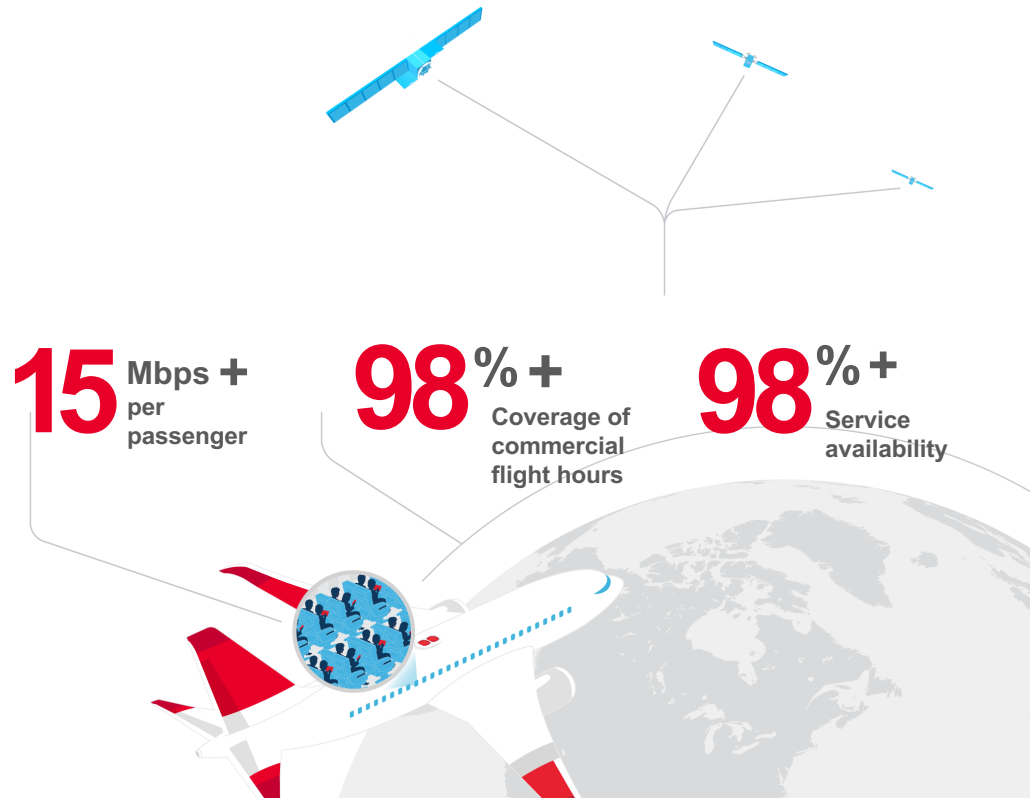


Gogo's Technology Roadmap

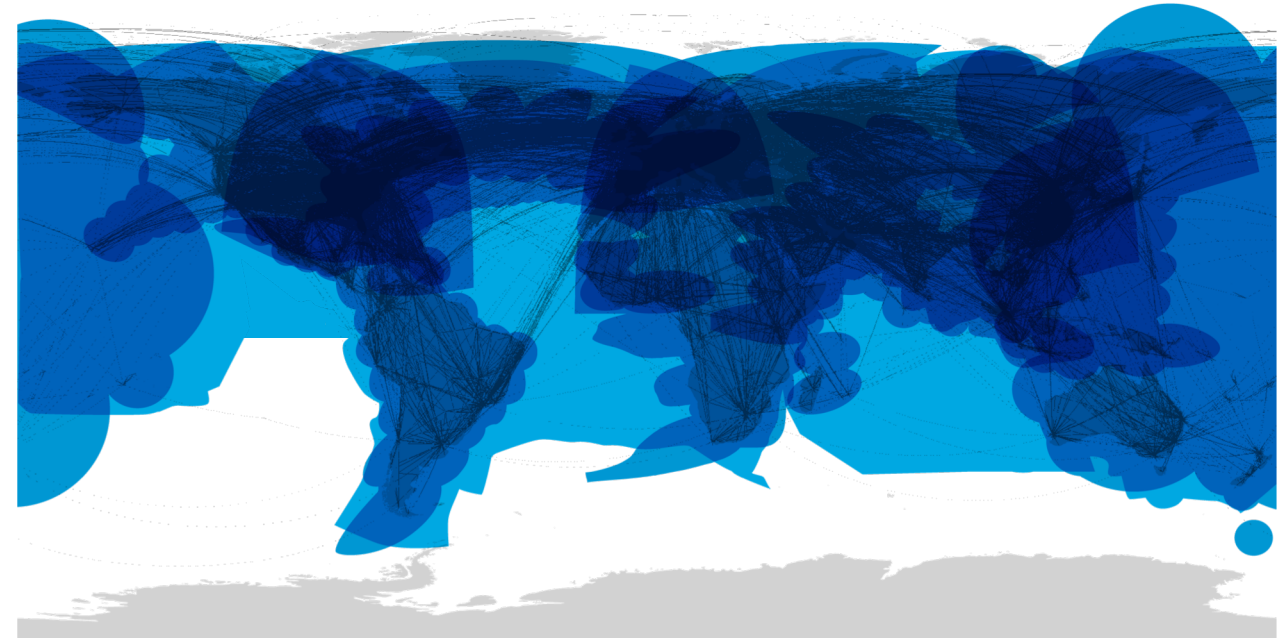


Performance everywhere with 2Ku

A ground-like experience



Everywhere aircraft fly



Backstory - Unified Data Platform (UDP)



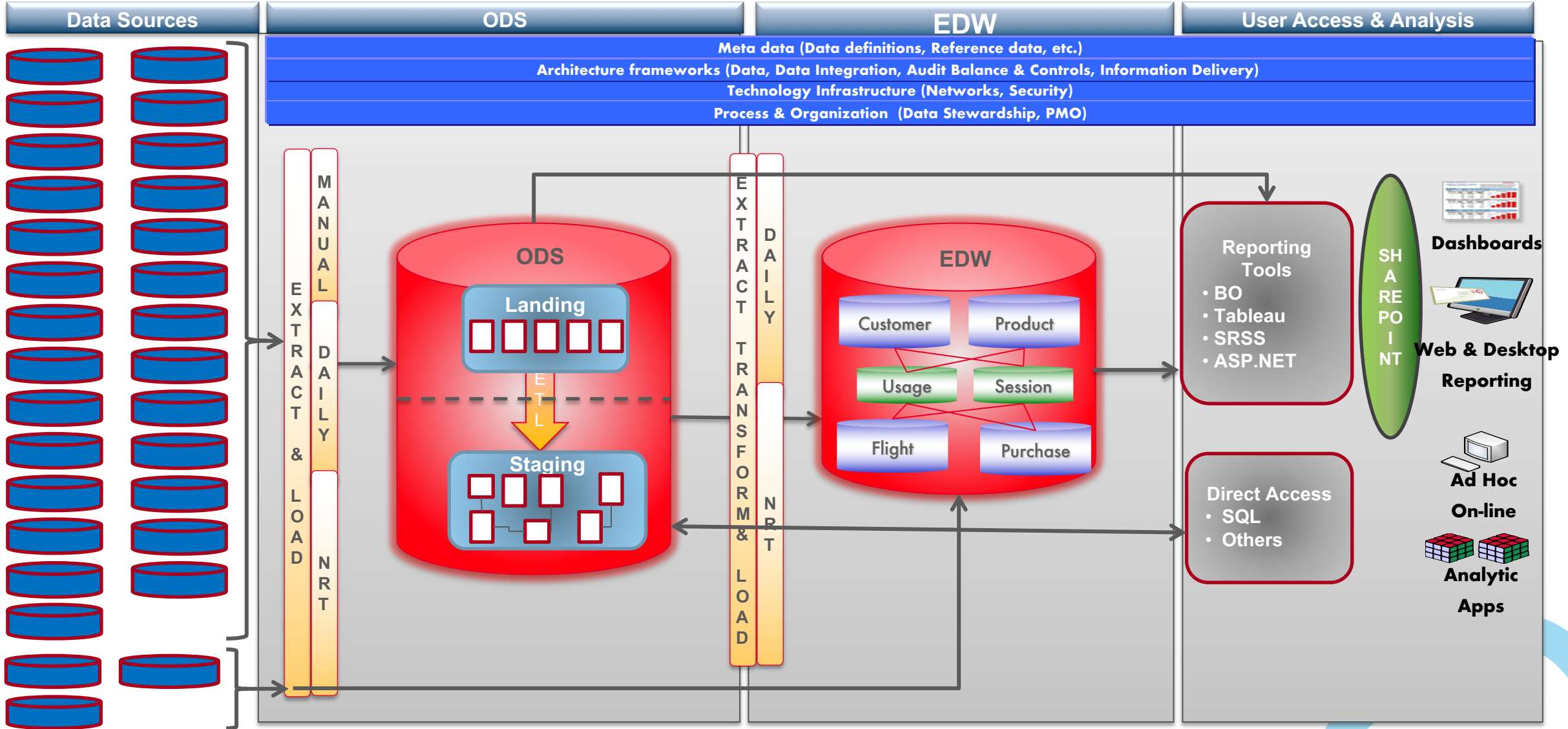
Goals:

- Event-driven data processing
- Near real-time analytics
- Prediction modeling

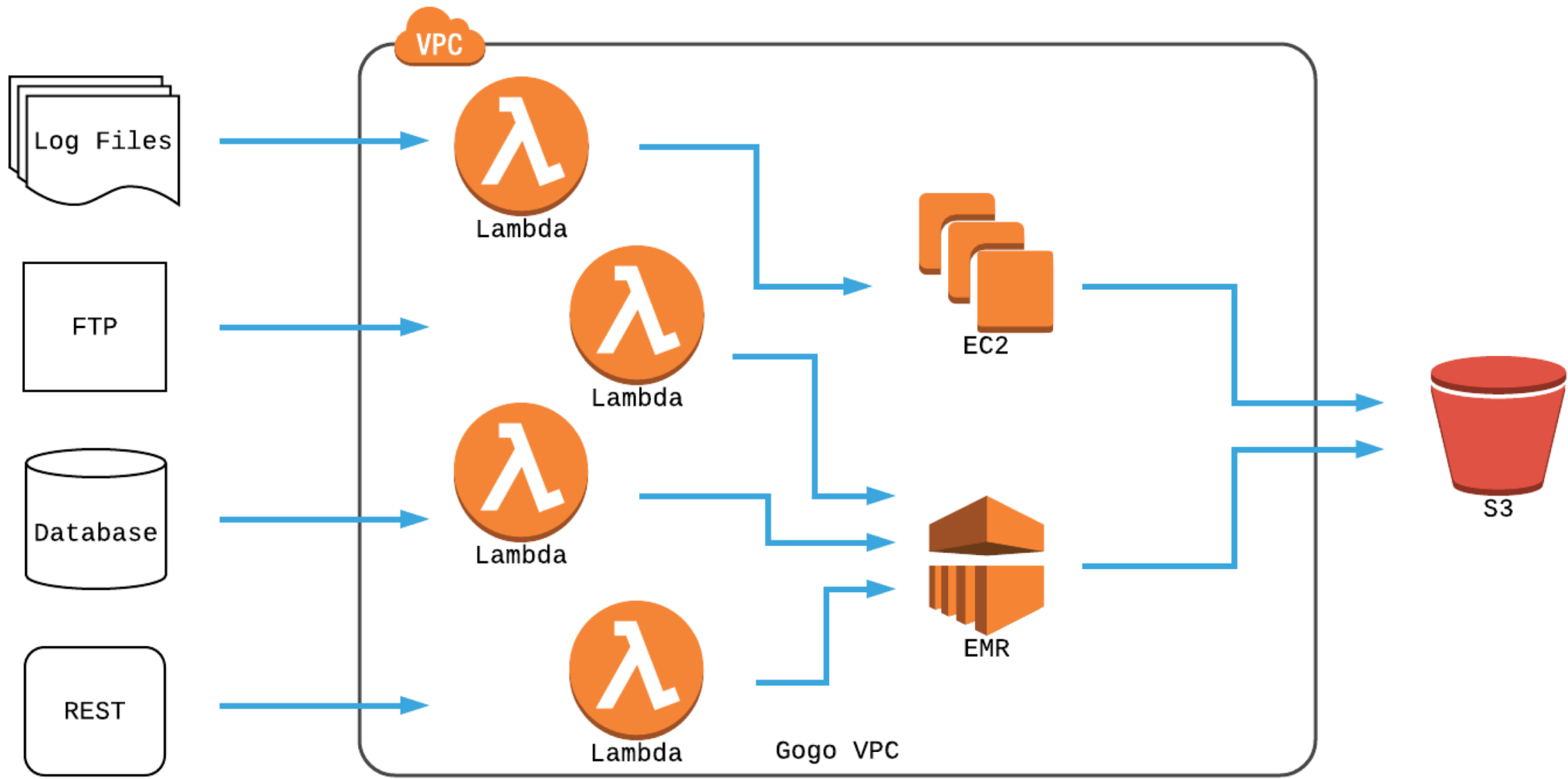
AWS:

- Availability
- Scalability
- Managed Services

Old world



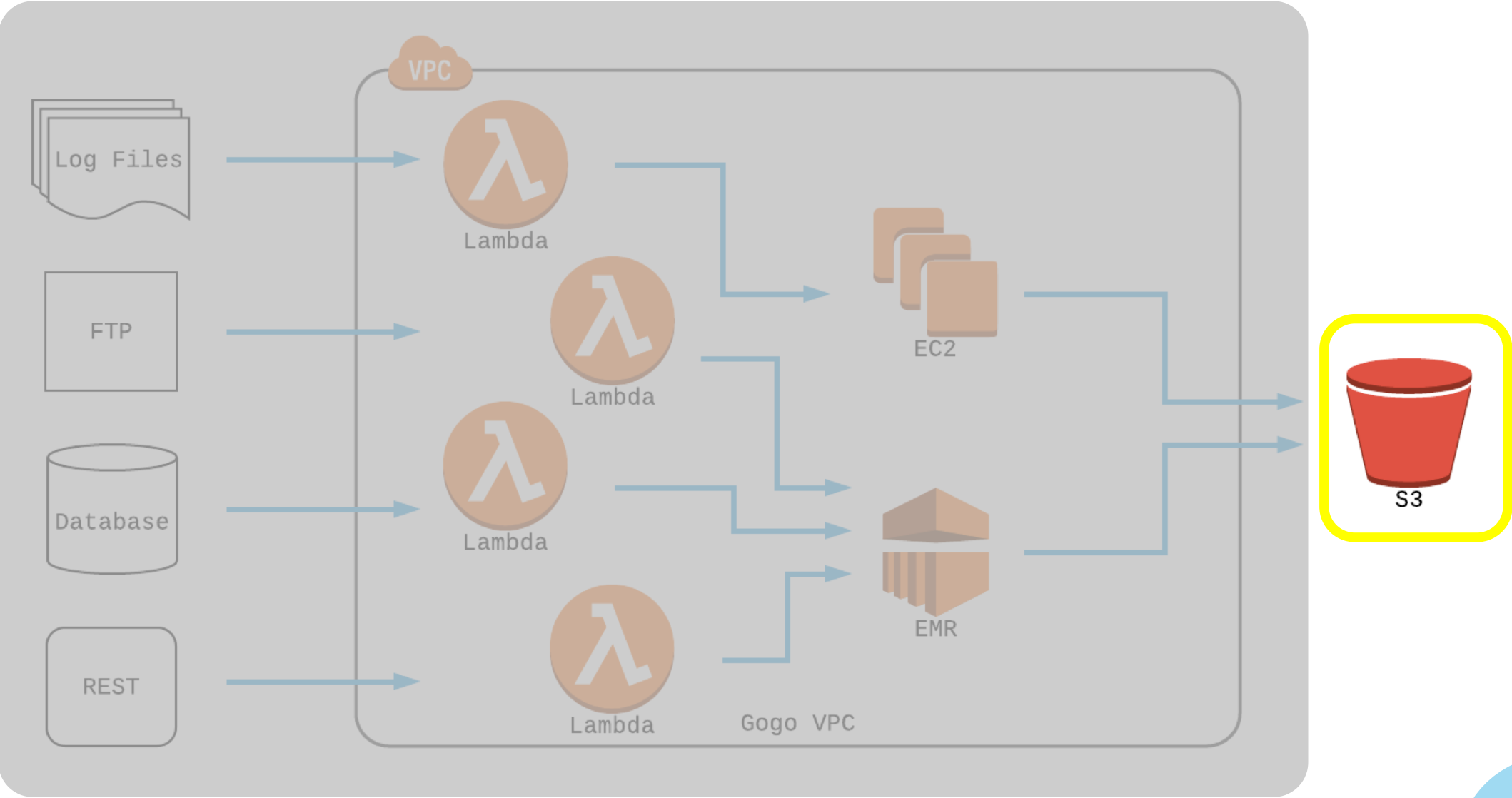
New world



Lessons Learned



S3: The object store



- S3 is an object store, not a file system
- Any “directory” structure is nothing more than a string prefix to the object

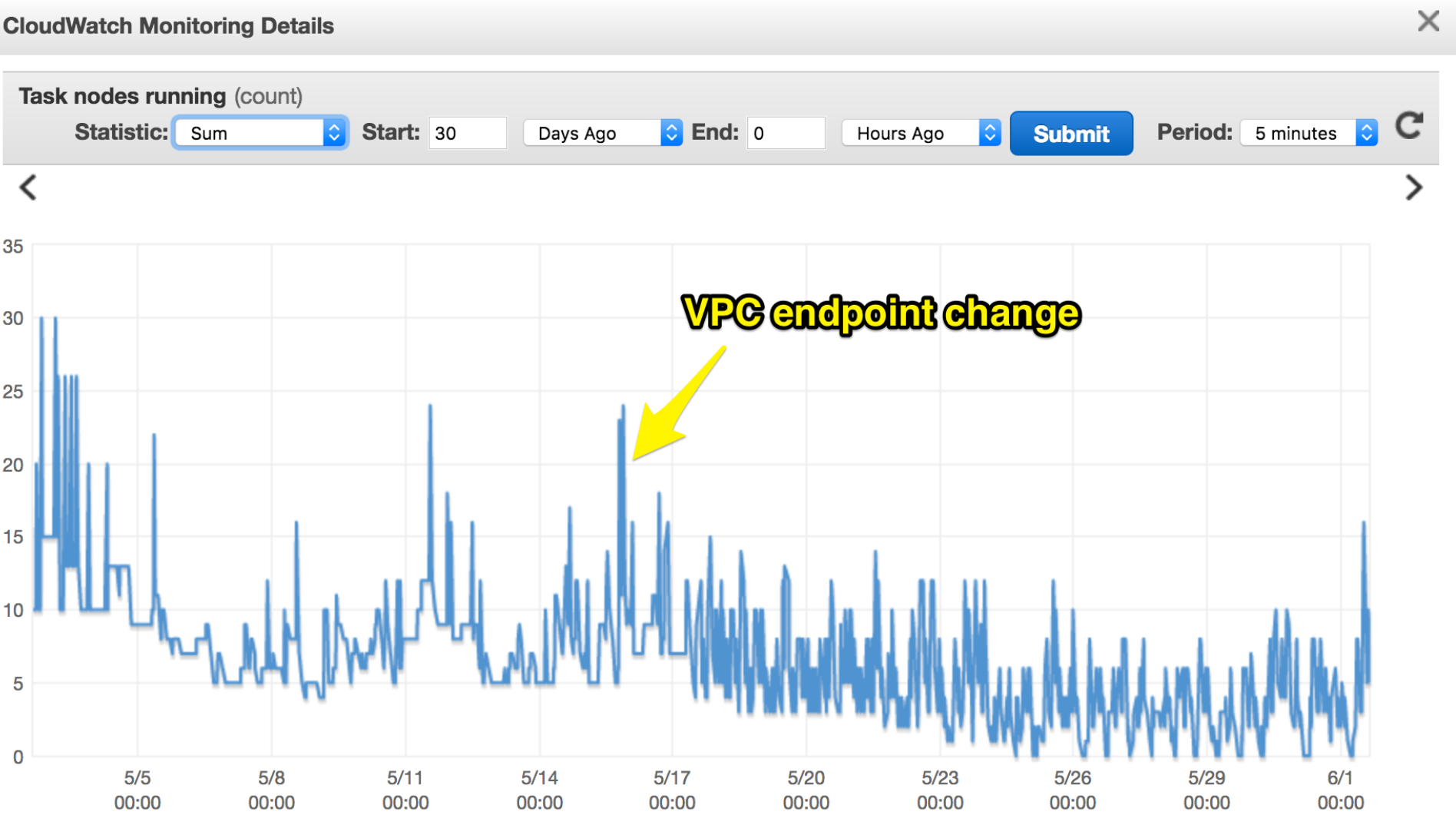
Issues:

- Bucket contention
- API Rate limiting (500 errors)

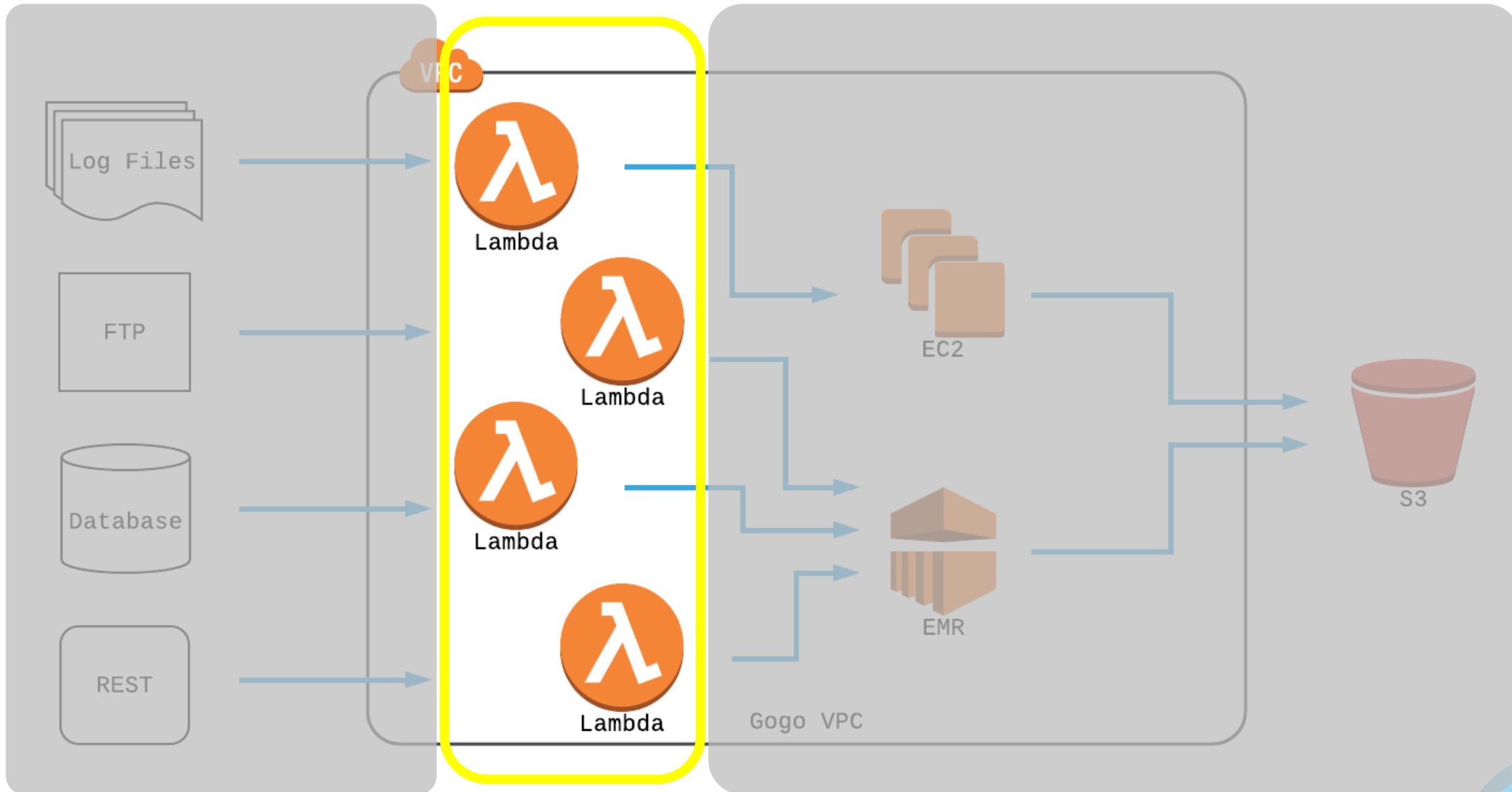
Parquet: The joys & the pitfalls

- Compressed, columnar file format
- Natively supported by many analytical tools
 - Pandas
 - Spark
 - Redshift
 - Presto/Drill/Impala
- Cost: High overhead when writing files

VPC Configuration



Lambda: Serverless computing at scale



Lambda: Pros and Cons, Caveats

Pros:

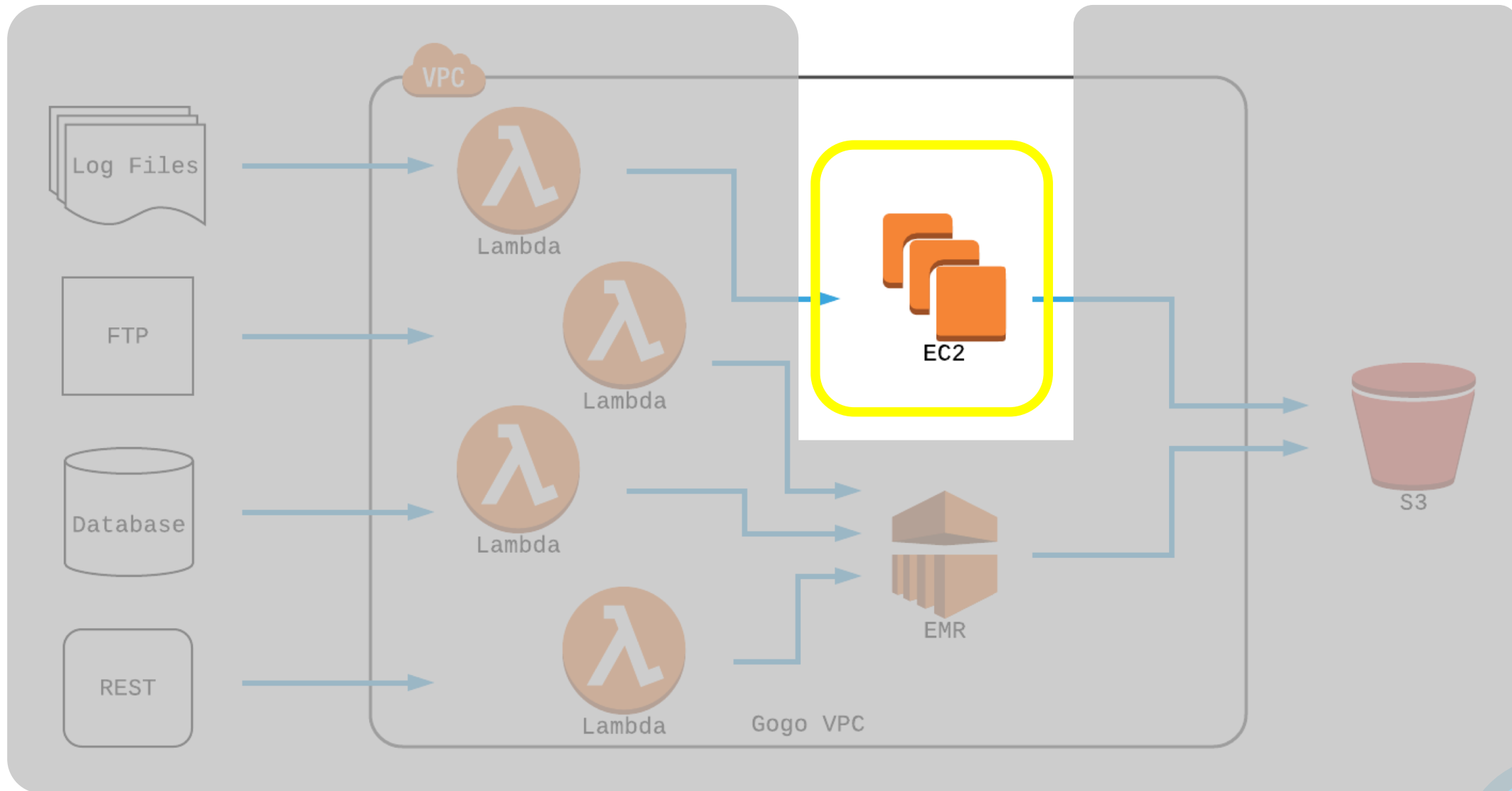
- No need to manage servers
- Integration with key AWS components (S3, SNS)
- Scalable

Cons:

- Difficult to test
- Environment is black box
- Scalable

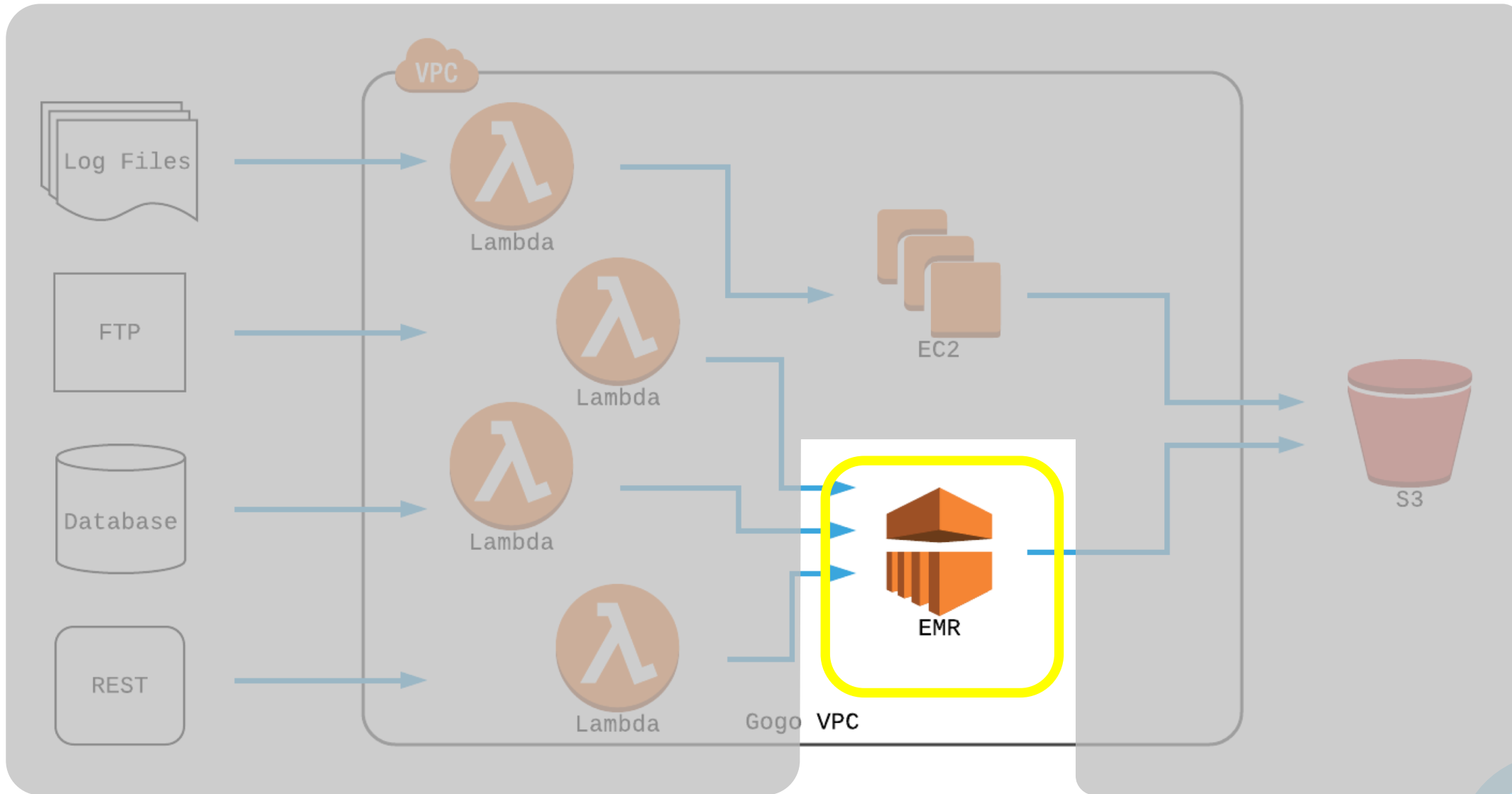
Caveats

- Cost
- Downstream processes (file volumes, API limits)
- Account limits/concurrency



- Sizing is hard, use more than Cloudwatch to gauge system metrics like memory
- Use spot whenever possible but be aware of the risks
- Spot-pricing can vary dramatically depending on AZ/Region
- Reserving for 1-year increments is a good balance of price/flexibility

EMR: Hadoop as a service



- Step API is limited, geared more toward short-lived clusters – Livy can help
- CI/CD is difficult
- Bigger might be more efficient (because of software costs)
- Scaling is slow/error-prone

Results

- Faster time to reports and analytics
- More ad-hoc analysis
- Similar cost

Other takeaways

- Take small bites
- Set a budget
- Create a local development framework (Docker)

Future State

- Containerization, microservices
- Move away from managed services
- Multi-cloud strategy (GCP, Azure)

Thank you

