# Scaling the Monolith

Daniel Larsen

Technical Evangelist

2019

# It's a Trilogy!

Scaling the monolith

Refactoring the monolith
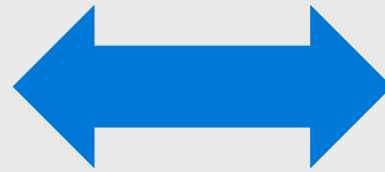
Strangling the monolith
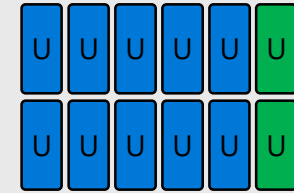
# Scalability

## Scale up

Scale up App Services
Scale up VMs
Scale up SQL DB
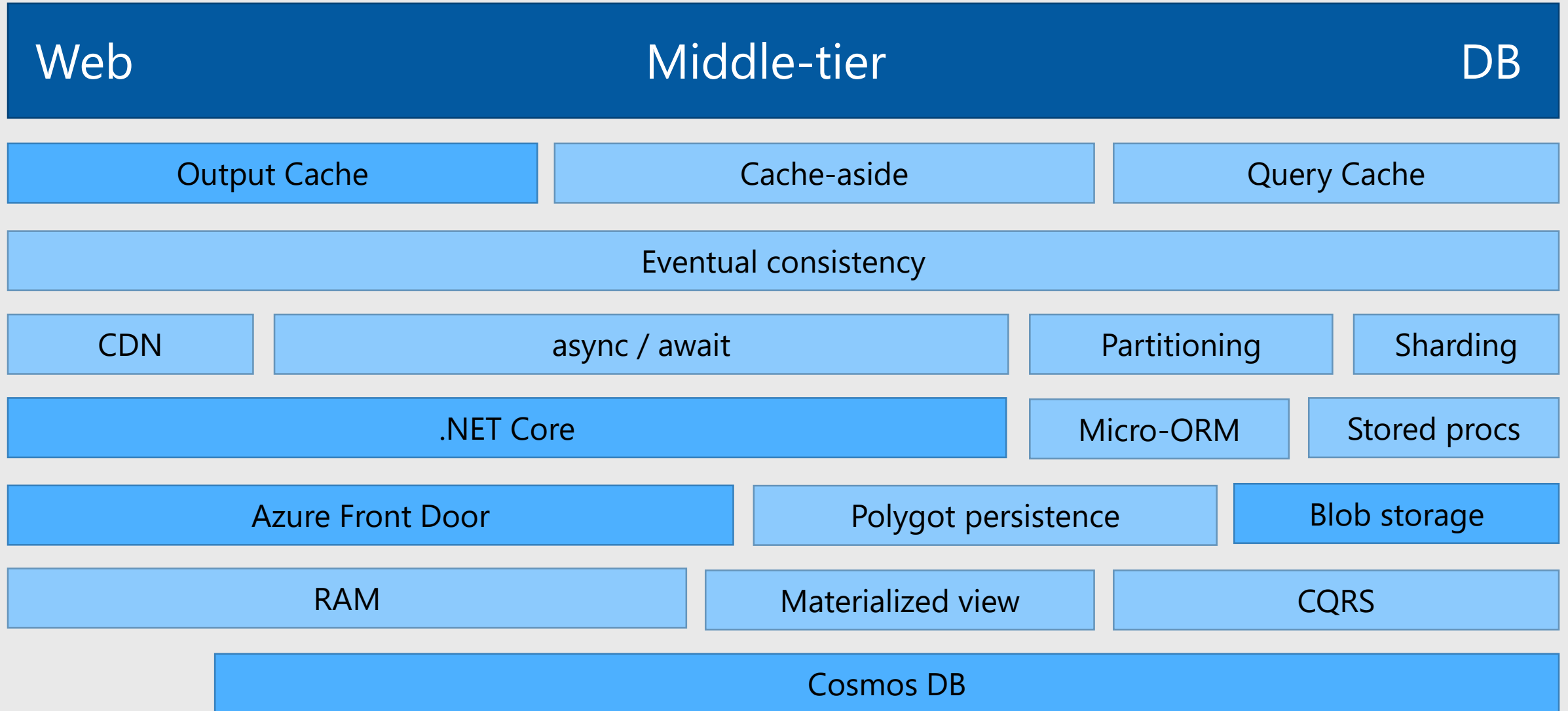More cores, more RAM

## Scale out

Scale out App Services
Scale out VM Scale Sets
Auto-scale
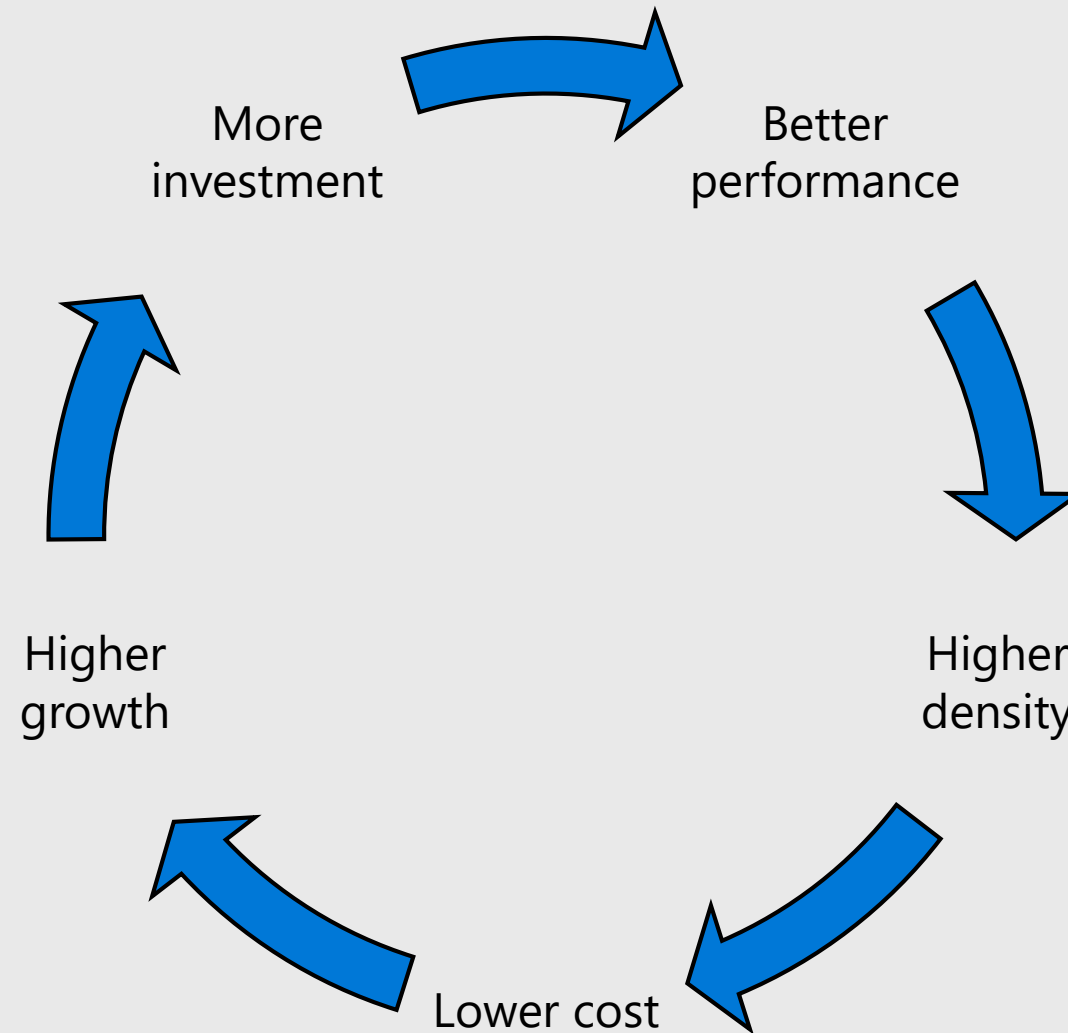Web farm + Load balancing
SQL Clustering

## Optimize

Scalability of Software Architecture
Optimization of Software (code)
Optimization of Platform services
Cloud design patterns

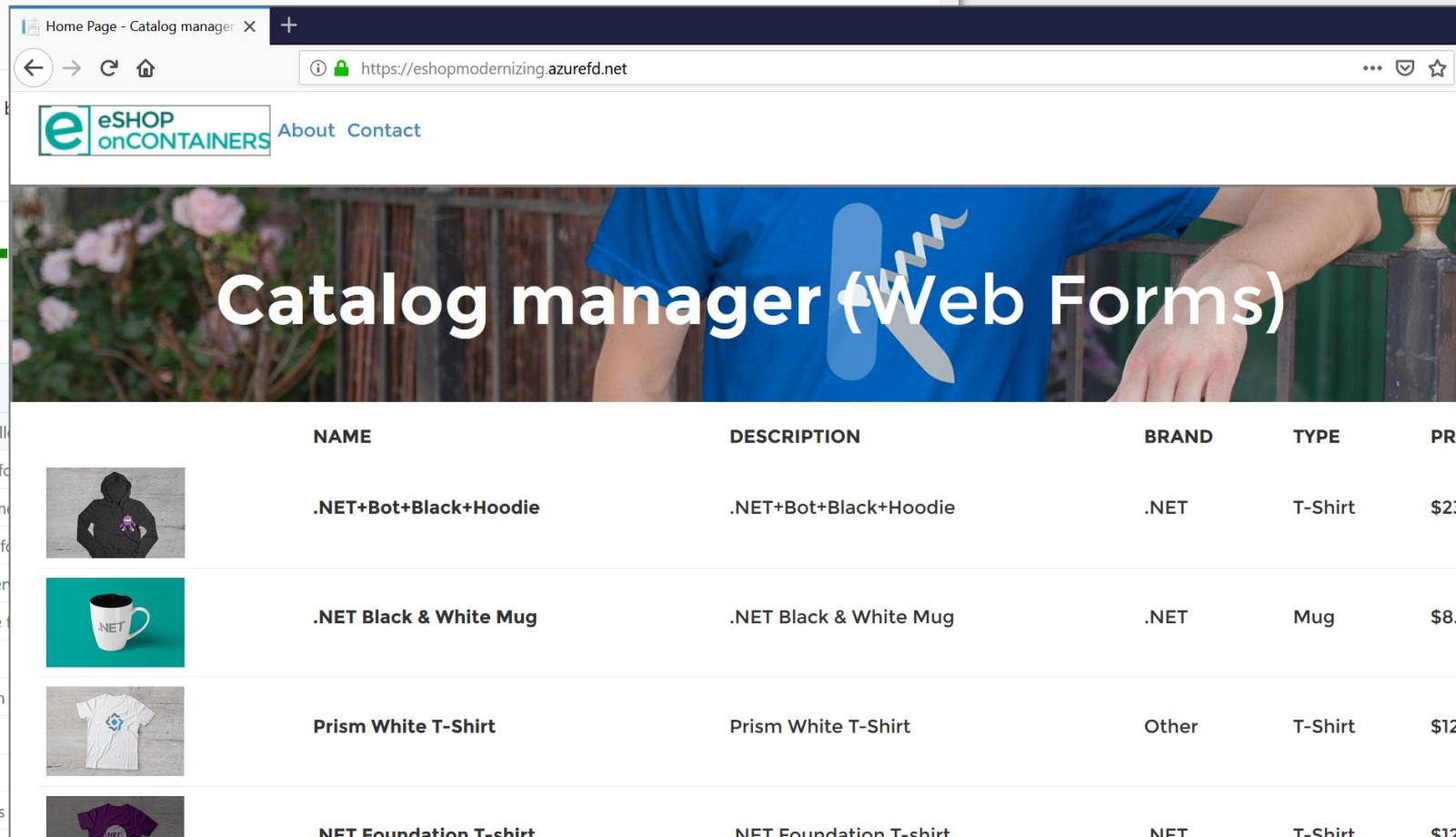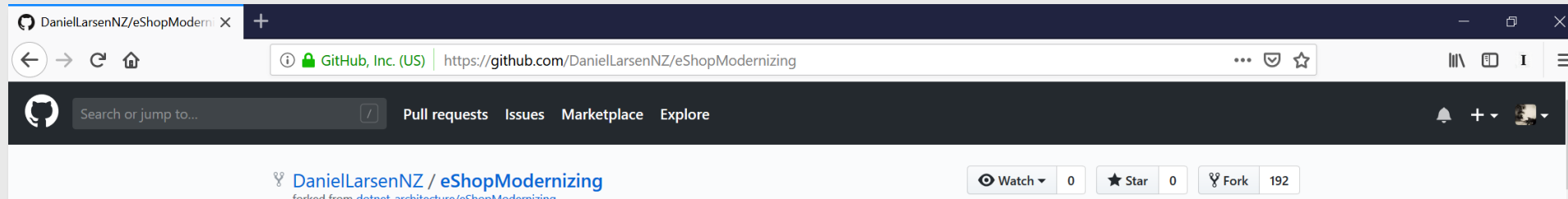# Patterns and practices for Scalability of software

| Web | Middle-tier | DB |
|-----|-------------|-----|

| Output Cache | Cache-aside | Query Cache |

| Eventual consistency |

| CDN | async / await | Partitioning | Sharding |

| .NET Core | Micro-ORM | Stored procs |

| Azure Front Door | Polygot persistence | Blob storage |

| RAM | Materialized view | CQRS |

| Cosmos DB |

# Performance engineering

More investment → Better performance → Higher density → Lower cost → Higher growth → More investment

# eShopModernizing

# Key performance indicators

↑10,000
rpm / core

↓10ms
server response

↓1s
page load

↓100
dtu max

↑25:1
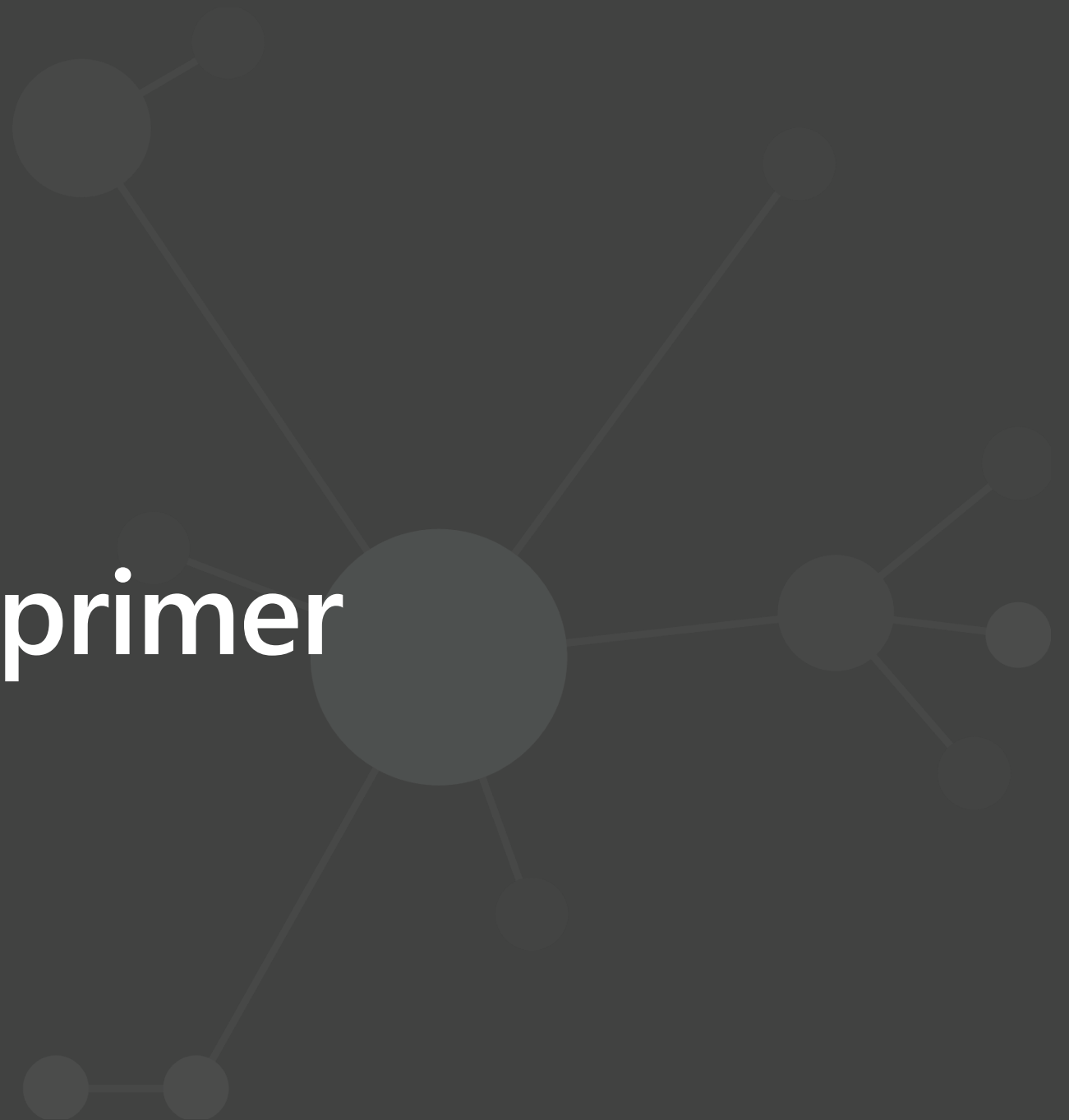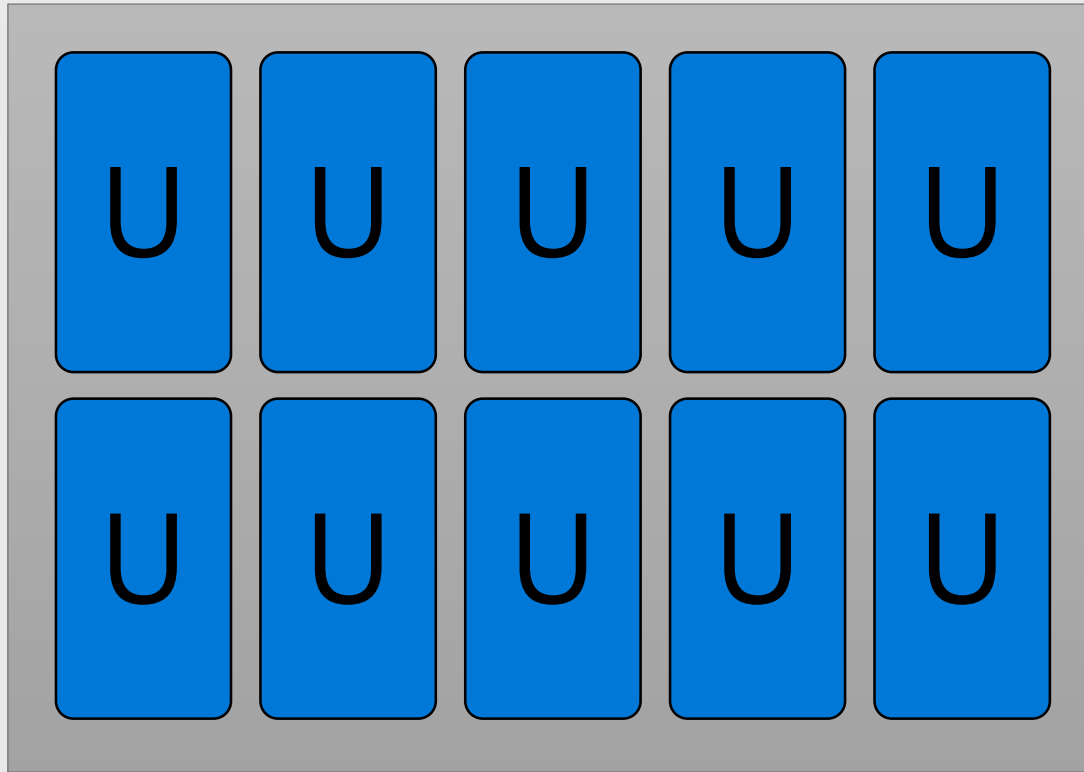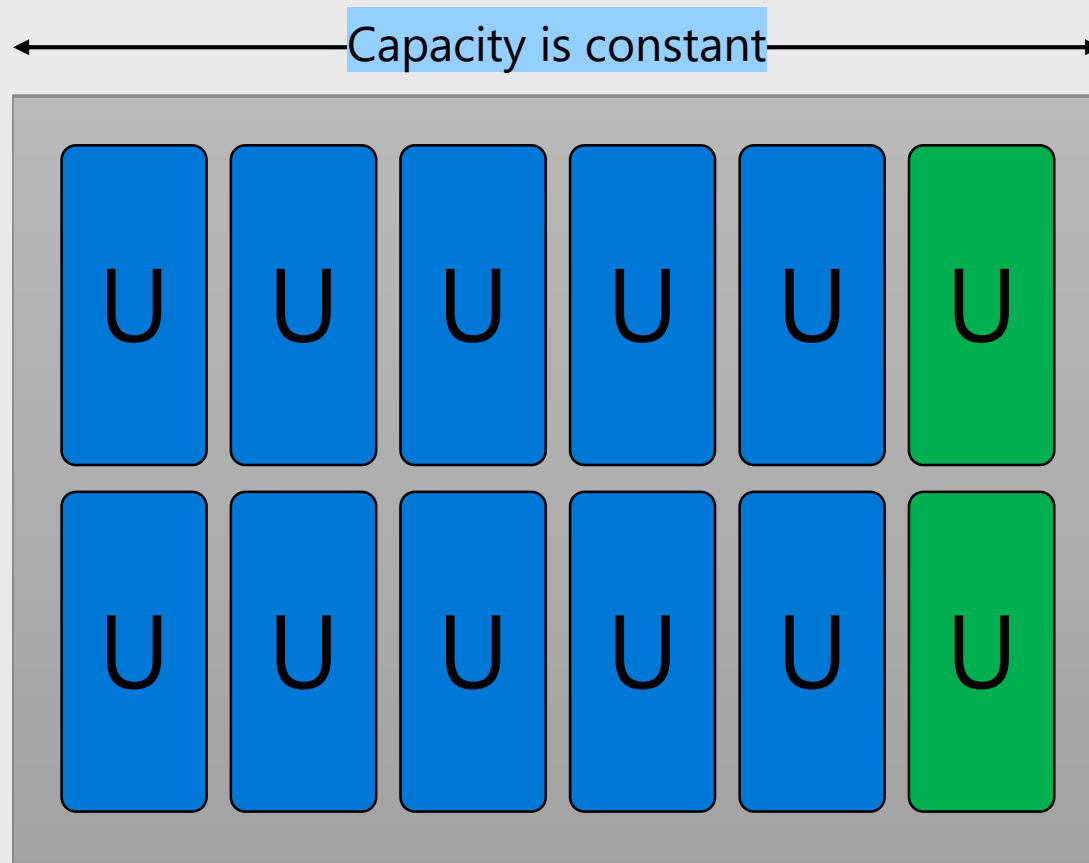web : db calls

# Cloud economics primer

# Density

Cost to serve: e.g. Cost per *n* users / tenants / transactions

# Density

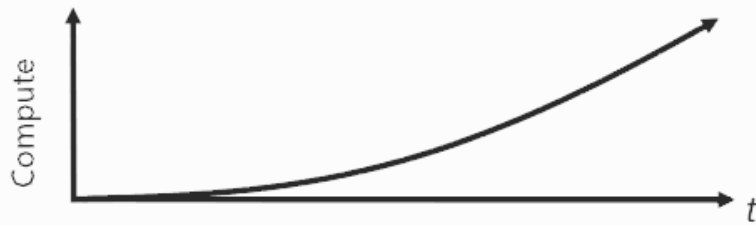**Cost to serve: e.g. Cost per *n* users / tenants / transactions**

Capacity is constant
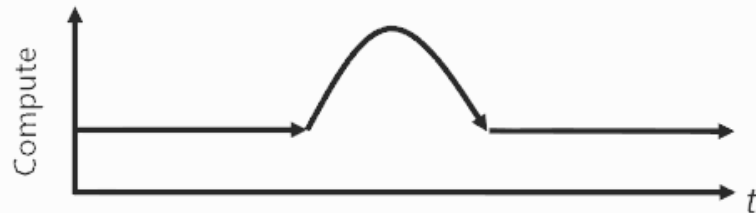
20% higher density = 20% reduction in cost to serve

# On and Off

On & off workloads (e.g. batch job)
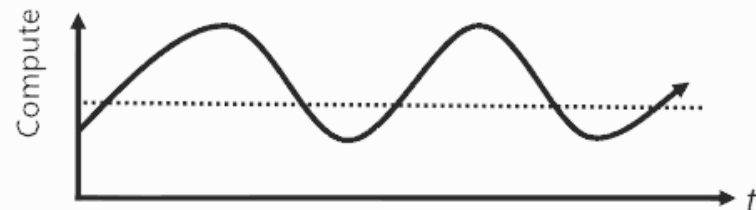Over provisioned capacity is wasted
Time to market can be cumbersome

# Growing Fast

Successful services needs to grow/scale
Keeping up w/ growth is big IT challenge
Cannot provision hardware fast enough

# Unpredictable Bursting

Unexpected/unplanned peak in demand
Sudden spike impacts performance
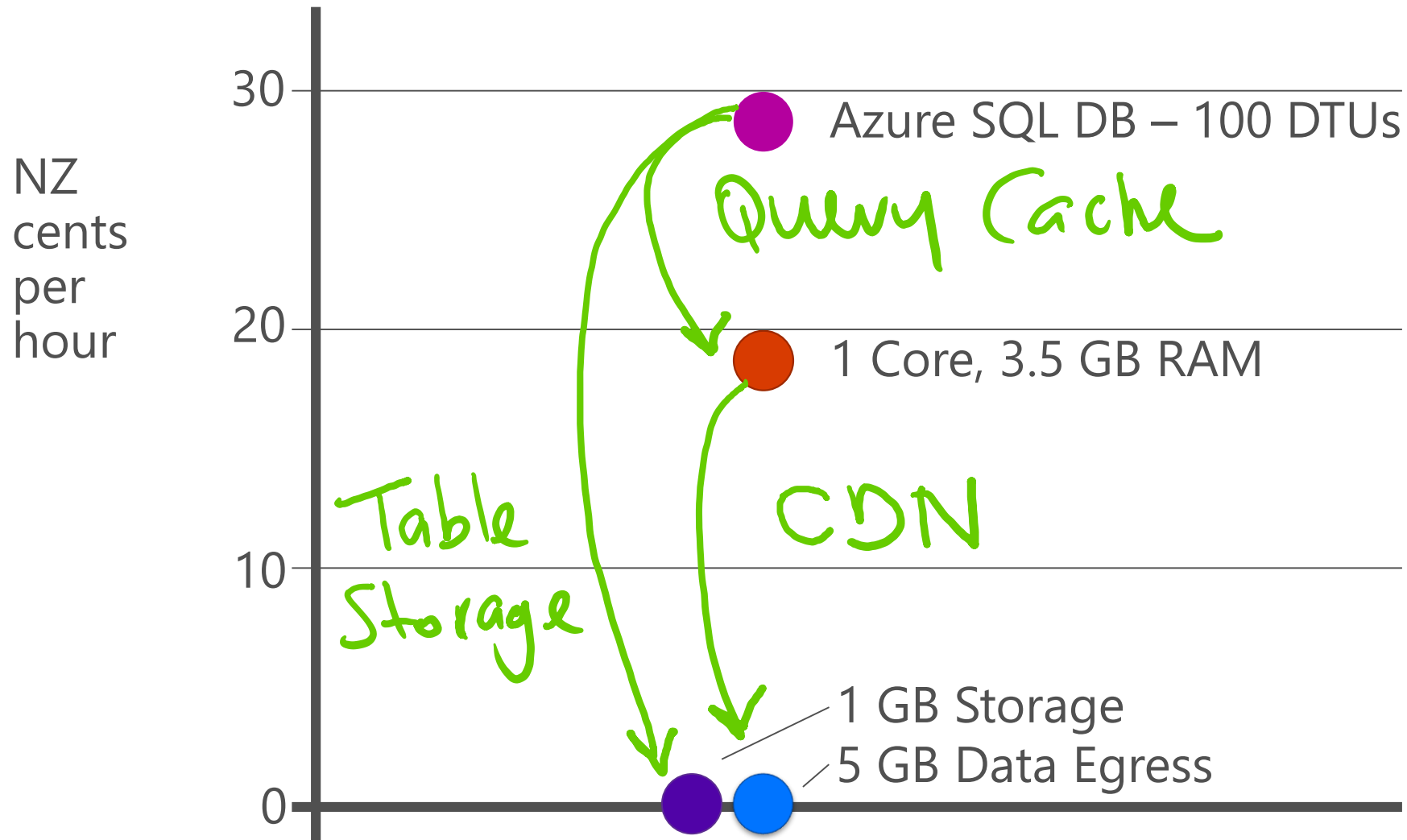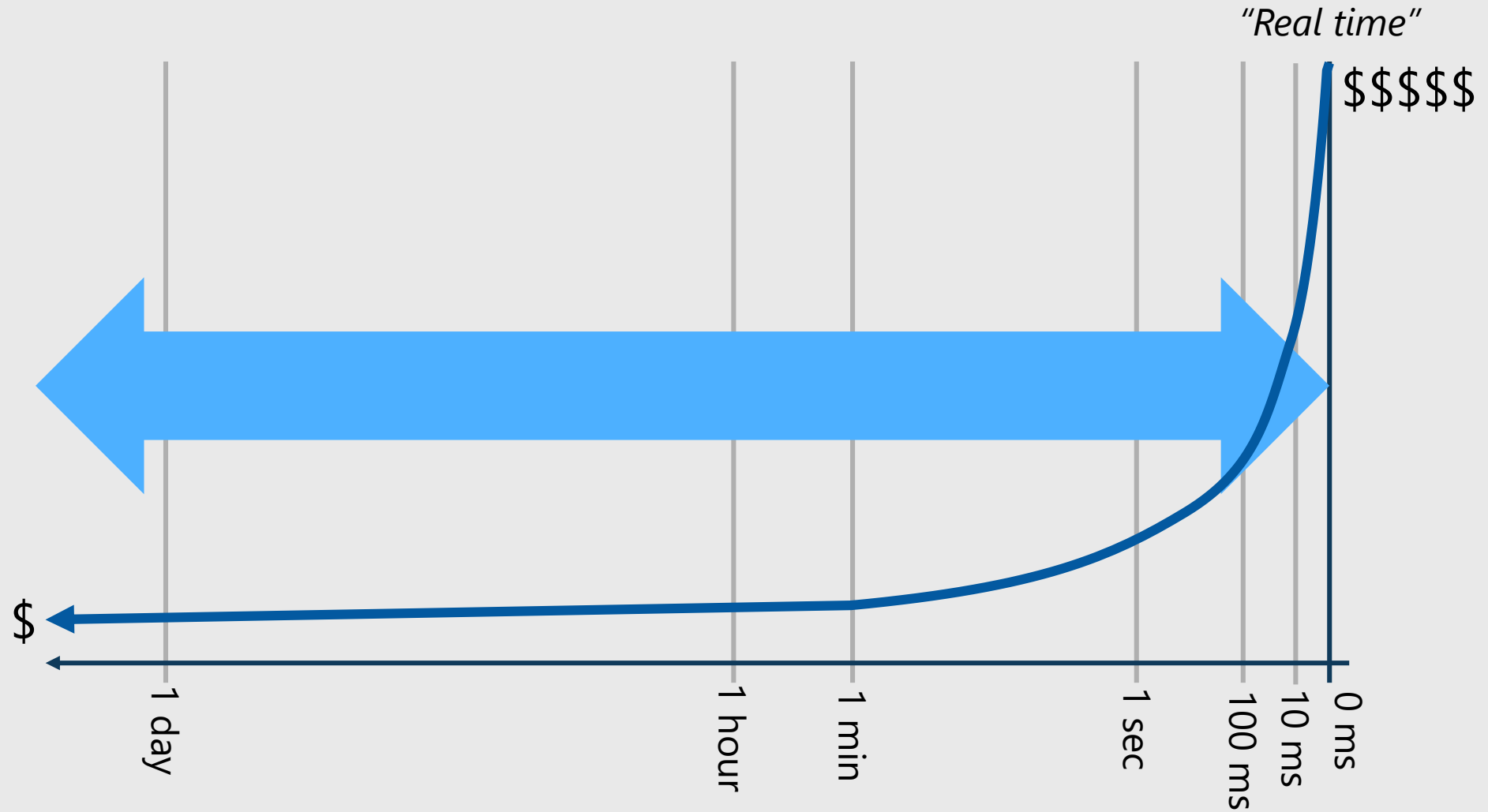Can't over provision for extreme cases

# Predictable Bursting

Services with micro seasonality trends
Peaks due to periodic increased demand
IT complexity and wasted capacity

# Cloud Computing Patterns

# Optimising costs

# Eventual consistency

# As a User



dalars / MyFirstProject / Boards / Backlogs

Search

M

≡ MyFirstProject Team ⌄ ☆ ⅄

📖 USER STORY 541*

541 | **As a Store User I want to browse a list of Products**

👤 Unassigned                           💬 0 comments   Add tag          💾 Save & Close ⌄   👁 Follow   ↻   ↩   ⋯   ✕

| State | ⚪ New | Area | MyFirstProject | | Updated just now |
|---|---|---|---|---|---|
| Reason | 🔒 New | Iteration | MyFirstProject\Release CHG1\Sprint 2 | | |

Details  🕐  🔗  📎

## Description

So that I can choose a product to add to my cart.

## Acceptance Criteria

1. List is paged in pages of 10 products (by default)
2. Displays pic for each Product
3. Product pic and Name click through to Product detail

## Non-functional requirements

1. TTL = 30 seconds
2. Avg. page load < 1.5 seconds

**B** *I* <u>U</u> ☰ ☰ 🖍⌄ A⌄ 🙂 ⇤ ⇥ s̶ AA⌄ </> A̶ ⋯

## Discussion

### Planning

Story Points

Priority
2

Risk

### Classification

Value area
Business

### Development

➕ Add link
Development hasn't started on this item.
Create a new branch

### Related Work

➕ Add link ⌄
There are no links in this group
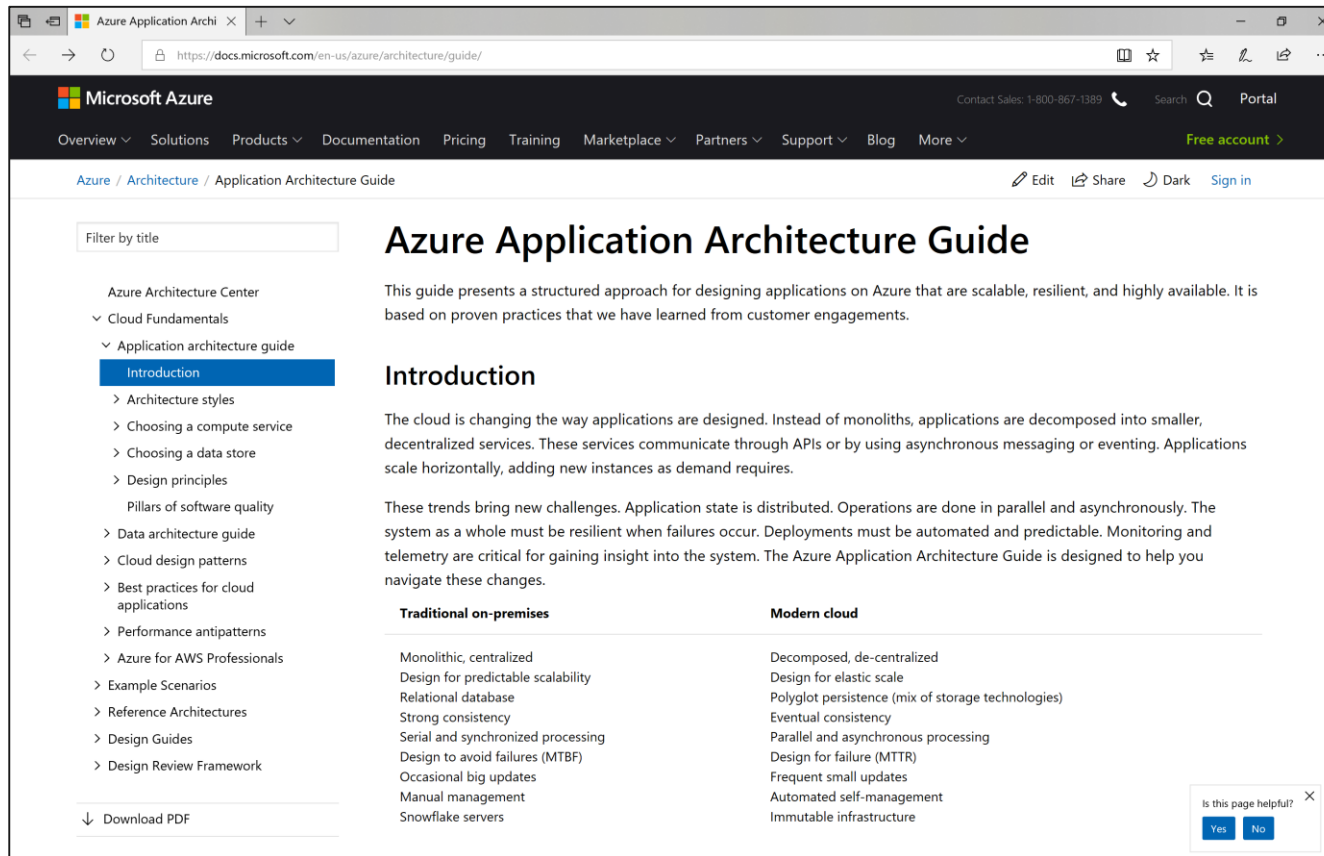
# Output Cache

# Scalability checklist

# Azure Application Architecture Guide



https://aka.ms/architecture

# Scalability checklist

The **Scalability checklist** contains 30+ recommendations. 10 are listed here.

| | Recommendation | Patterns & guidance |
|---|---|---|
| 1 | Review the performance antipatterns | [Performance antipatterns](#) |
| 2 | Use asynchronous calls | [Asynchronous programming](#) |
| 3 | Design for eventual consistency | [Data Consistency Primer](#) |
| 4 | Use data partitioning, Consider de-normalizing data | [Data partitioning](#), [Materialized View](#) |
| 5 | Minimize load on the data store, Minimize the volume of data retrieved, Optimize and tune SQL queries and indexes | [Automatic tuning in Azure SQL Database](#) |
| 6 | Partition the workload, design for scaling, scale as a unit | [Microservices architecture](#) |
| 7 | Aggressively cache, Use output cache, Enable client caching | [Caching Guidance](#) |
| 8 | Offload and distribute intensive CPU/IO tasks | [Background jobs](#), [Competing Consumers](#) |
| 9 | Use queues to level load | [Queue-Based Load Leveling Pattern](#) |
| 10 | Carry out performance profiling and load testing | [Testing cloud service performance](#) |

https://docs.microsoft.com/en-us/azure/architecture/checklist/scalability

# Performance antipatterns for cloud applications

A *performance antipattern* is a common practice that is likely to cause scalability problems when an application is under pressure.

| | |
|---|---|
| Busy Database | Offloading too much processing to a data store |
| Busy Front End | Moving resource-intensive tasks onto background threads |
| Chatty I/O | The cumulative effect of a large number of I/O requests |
| Extraneous Fetching | Retrieving more data than is needed, resulting in unnecessary I/O |
| Improper Instantiation | Repeatedly creating and destroying objects that are designed to be reused |
| Monolithic Persistence | Using the same data store for data with very different usage patterns |
| No Caching | Failing to cache data |
| Synchronous I/O | Blocking the calling thread while I/O completes |

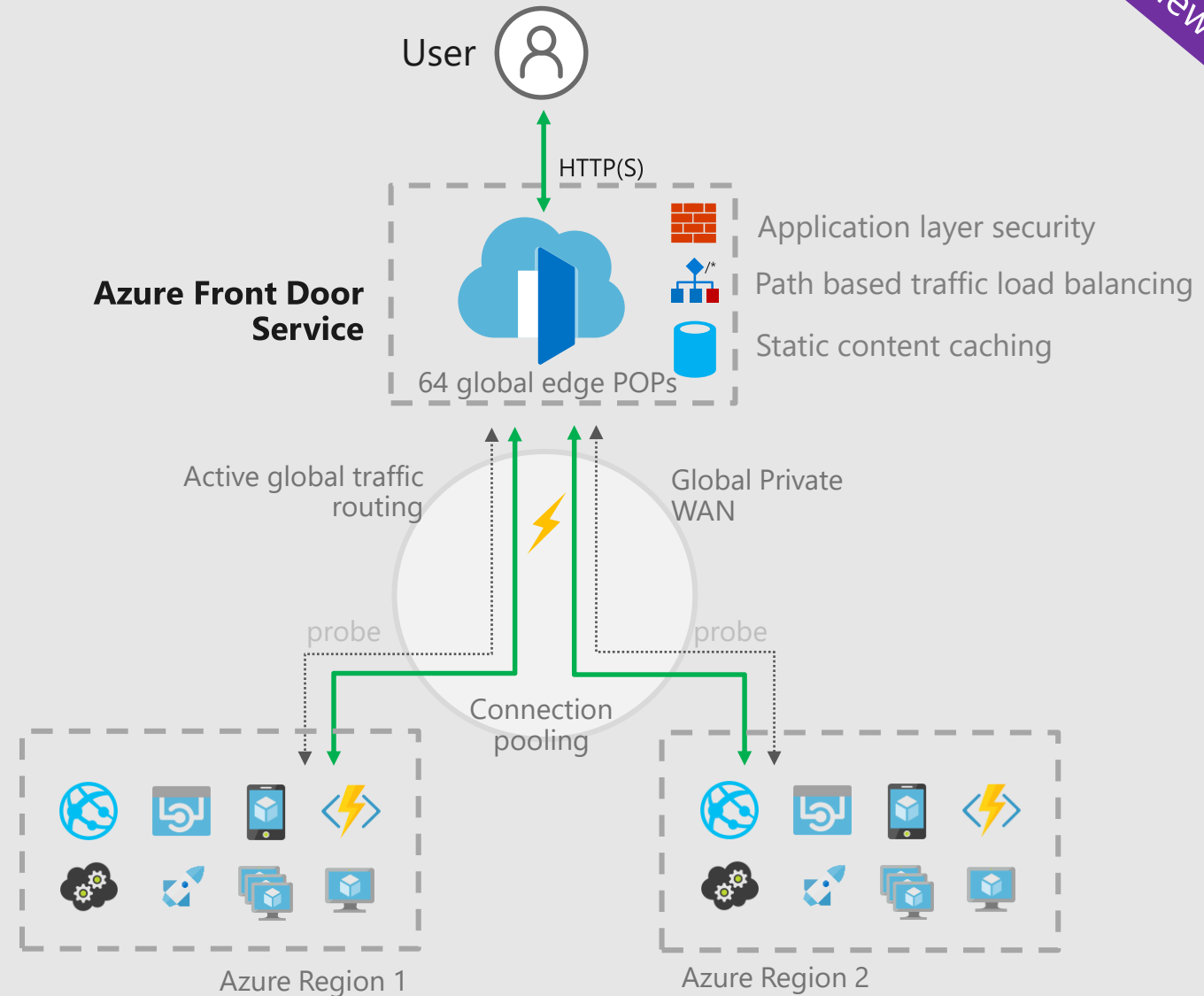https://docs.microsoft.com/en-us/azure/architecture/antipatterns/index

# Azure Front Door Service

**Your secure entry point for delivering globally performant hyperscale apps.**

✓ Application acceleration at Microsoft's edge

✓ Integration with App Services

✓ Global HTTP load balancing with instant failover

✓ Massive SSL offload

✓ Integrated static content caching

✓ Central application traffic dashboard

Office 365  Azure  Skype  Bing
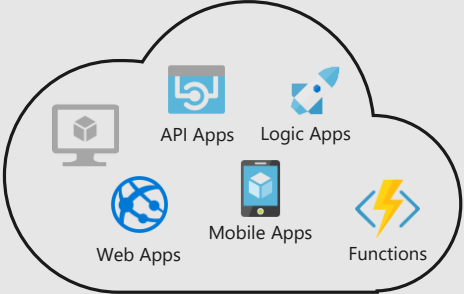Azure DevOps  MSN OneDrive
Xbox Cortana Windows Teams

Microsoft



Preview

User

HTTP(S)

**Azure Front Door Service**

64 global edge POPs

Application layer security

Path based traffic load balancing

Static content caching

Active global traffic routing

Global Private WAN

probe

probe

Connection pooling

Azure Region 1

Azure Region 2

https://azure.com/frontdoor

# Connection establishment and response

API Apps
Logic Apps
Web Apps
Mobile Apps
Functions

Application backend

Clients

syn

syn.ack

ack

TCP Connect

Client Hello

Server Hello, Certificate, Server Hello Done

SSL Connect

Client Key Exchange, Cipher Spec, Client Finished

Change Cipher Spec, Finished

Time to First Byte

GET

App Latency

First Byte

Total Response Time

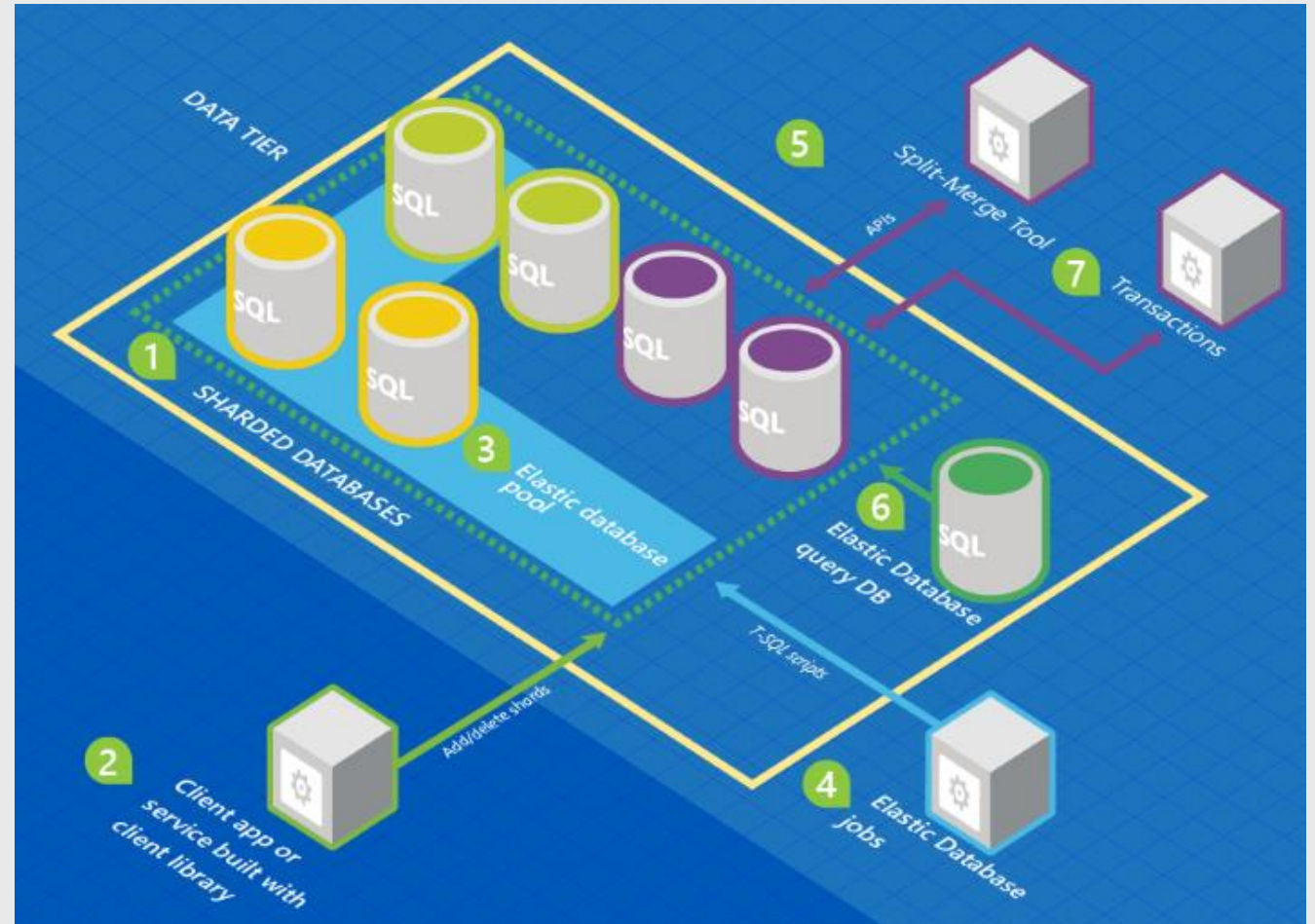ack

ack

Download Time

ack

Demo
# Azure Front Door

# Scaling out with Azure SQL Database
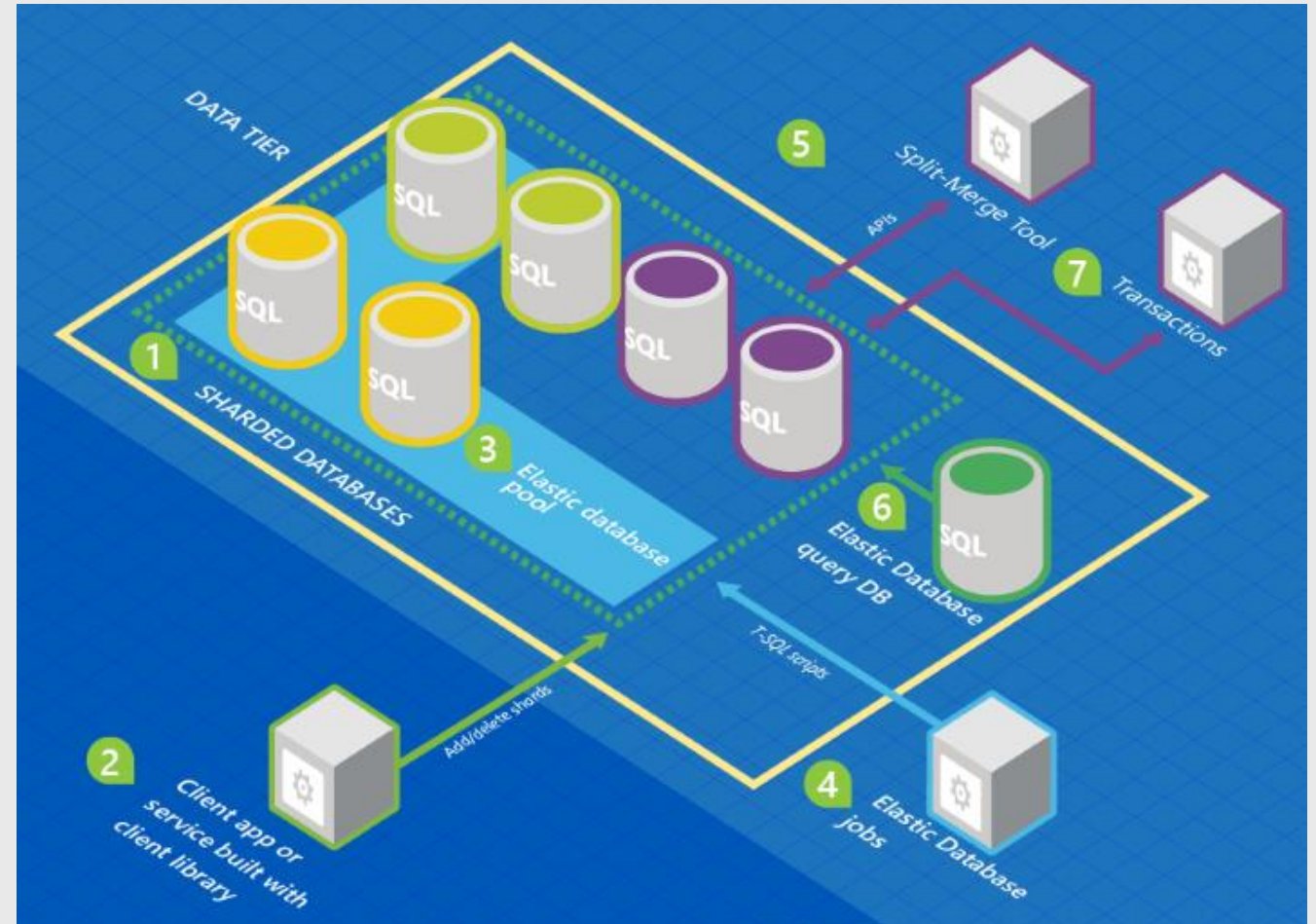## Elastic Database Tools

1. Azure SQL DBs – Shard set

2. Elastic Database client library

3. Subset of DBs can be put into an Elastic pool

4. Elastic Database job runs scheduled or ad hoc T-SQL scripts against all databases

5. Split-merge tool to move data from one shard to another

6. Elastic Database query to write a query that spans all databases in the shard set

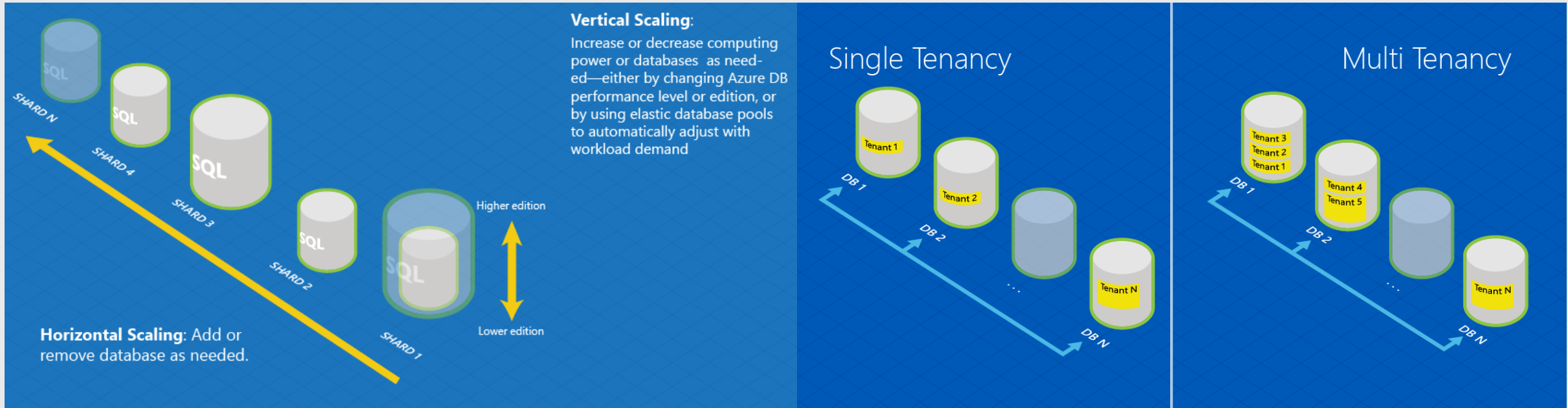7. Elastic transactions to run transactions that span several databases

# Scaling out with Azure SQL Database
## Elastic Database Tools

1. Azure SQL DBs – Shard set

2. Elastic Database client library

3. Subset of DBs can be put into an Elastic pool

4. Elastic Database job runs scheduled or ad hoc T-SQL scripts against all databases

5. Split-merge tool to move data from one shard to another

6. Elastic Database query to write a query that spans all databases in the shard set

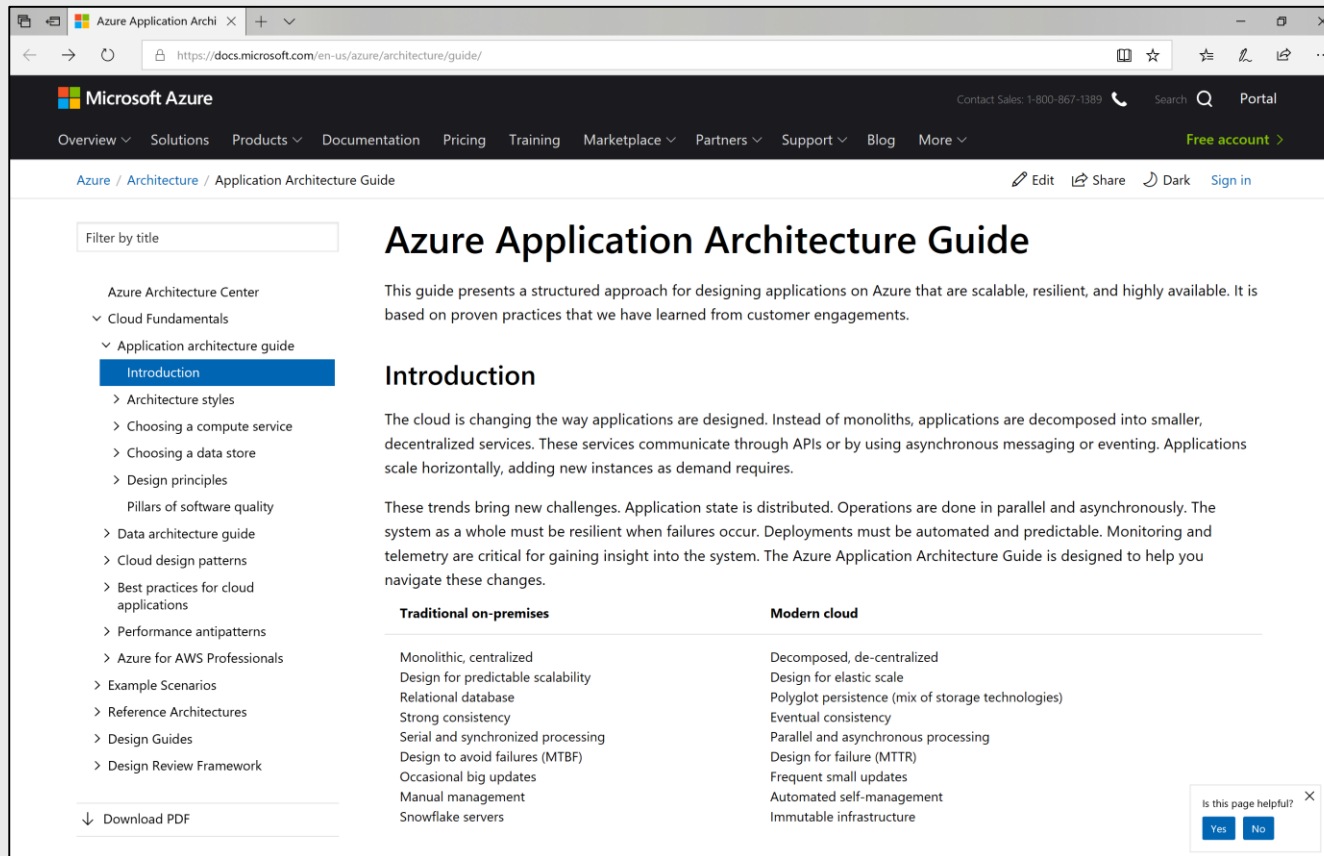7. Elastic transactions to run transactions that span several databases

# Sharding



Vertical Scaling: Increase or decrease computing power or databases as needed—either by changing Azure DB performance level or edition, or by using elastic database pools to automatically adjust with workload demand

Horizontal Scaling: Add or remove database as needed.

Single Tenancy

Multi Tenancy

https://docs.microsoft.com/en-us/azure/sql-database/sql-database-elastic-scale-introduction

# Wrapping up

# Azure Application Architecture Guide



https://aka.ms/architecture

*Thank you!*

Microsoft

https://github.com/DanielLarsenNZ/talks

# Image credits

https://commons.wikimedia.org/wiki/File:Glorious-blue-mountain-range.jpg