

# Well-Architected Cost Pillar

## Executive Watch-outs and Fixes

Cost guidance is great - until side effects eat the savings. This pocket-book explains each common cost pillar recommendation in plain words, flags the risk to watch, and gives five fixes you can apply fast.

---

### How to use

- In each review, find the recommendation you are applying.
- Read **What it means** to align on intent.
- Read the **Risk** so the failure mode is explicit.
- Pick one **Top mitigation** to implement now and one to stage next sprint.

## Use managed services by default

---

### What it means:

Prefer cloud-run databases, queues, runtimes, and analytics so teams focus on features. You accept the provider's upgrade cadence and defaults.

### Risk:

Forced upgrades or shifting defaults break compatibility, change quotas or performance, and raise cost without code changes.

### Top mitigations

#### Choose controlled tracks

Prefer LTS or manual upgrade channels, set maintenance windows, and pin versions in IaC and SDKs.

#### Test before it tests you

Run a canary on the next version, replay critical transactions, and smoke test backup and restore.

#### Stabilize interfaces

Add an adapter layer so service API changes do not ripple through your codebase.

#### Prepare safe rollback

Take snapshots, export configs, and keep a rehearse-ready rollback script for stateful services.

#### Escalate when needed

If cadence or compliance risk is unacceptable, run the service in-house behind the same interface with a migration plan.

## Buy commitments for steady load

---

### What it means:

Lower unit rates by committing to baseline usage (Savings Plans, Reserved Instances, Committed Use Discounts) for predictable workloads.

### Risk:

Overcommitment or wrong shapes strand discounts, block migrations, and hide falling utilization when traffic shifts.

### Top mitigations

#### **Baseline first**

Measure steady load for 30 to 60 days and size commitments to that baseline.

#### **Phase purchases**

Buy in tranches and expand only after utilization stays above target.

#### **Monitor like an SLO**

Track coverage and utilization monthly with alert thresholds and owners.

#### **Keep flexibility**

Prefer convertible or spend-based options and plan exchanges before shapes change.

#### **Tie to decommission plans**

Align commitments with lifecycle roadmaps so retirements do not strand discounts.

# Analyze and attribute expenditure

---

## What it means:

Tag and structure resources so spend maps cleanly to owners, products, and environments. Use accounts, projects, and subscriptions to mirror the organization.

## Risk:

Missing or inconsistent tags create orphaned spend, stall reviews, and cause finance vs. engineering disputes that block action.

## Top mitigations

### Enforce tagging at the gate

Require tags or labels in CI/CD and block deploys that miss the minimal key set.

### Mirror the business

Align hierarchy with product lines and environments for clean rollups.

### Backfill fast

Schedule jobs to repair missing tags on existing resources and log exceptions.

### Show owners their bill

Publish owner-level dashboards with the same 3 KPIs used in executive reviews.

### Make it policy

Document required keys, allowed values, and validation points in platform guardrails.

## Adopt a consumption model and autoscaling

---

### What it means:

Pay for what you use and scale automatically up and down. Aim for zero idle by setting sensible min and max capacity.

### Risk:

Mis-set scaling floors or new service behavior lock in 24x7 capacity, or quota limits force manual overrides that spike cost.

### Top mitigations

#### **Bound the range**

Declare min and max capacity in IaC and avoid high floors.

#### **Load test new behavior**

Test autoscaling and quotas under realistic traffic before rollout.

#### **Subscribe to limits**

Track provider notices for limit and scaling changes and record them in release notes.

#### **Protect with budgets**

Set cost and utilization alerts tied to scaling metrics.

#### **Ship safe defaults**

Make templates favor horizontal scaling and zero-idle patterns.

## Design for multi-region or hybrid

---

### What it means:

Place workloads across regions or clouds to meet latency, resilience, or residency needs.

### Risk:

Egress fees, duplicated storage, and chatty cross-region calls inflate cost and add complex failure modes.

### Top mitigations

#### **Map data flows early**

Include a simple data path diagram and egress estimate in design reviews.

#### **Keep compute with data**

Collocate services with primary datasets and avoid cross-region chatty patterns.

#### **Use caches and CDNs**

Serve repeats from edge or in-region caches to cut transfer volume.

#### **Control links**

Use private links or interconnects where appropriate and negotiate discounts for predictable flows.

#### **Set guardrails**

Budget and alert on egress per workload and block cross-region use that exceeds thresholds.