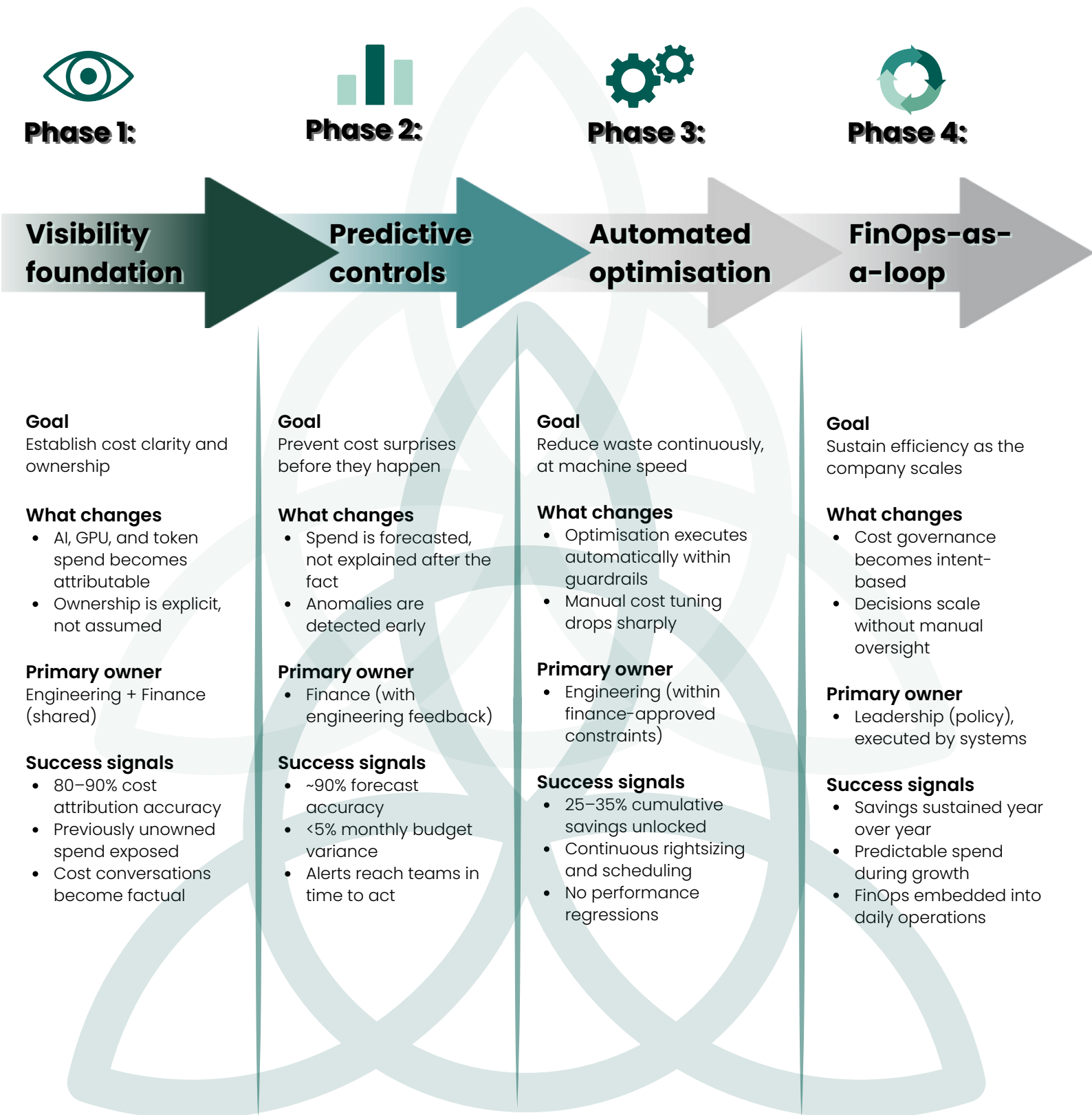


The 4-phase AI FinOps roadmap

From cost visibility to autonomous optimisation

How high-performing startups govern AI spend in 2026



Sequencing matters.

Teams that follow this progression consistently unlock 25–35% savings within 12 months.

AI FinOps: 60-day execution checklist

A week-by-week playbook designed for fast wins without slowing delivery

Phase 1:

weeks 1-2

Establish cost ownership

Phase 2:

weeks 3-4

Move from reporting to prediction

Phase 3:

weeks 5-6

Activate automated optimisation

Phase 4:

weeks 7-8

Close the leadership loop

Objective

Make AI and cloud spend visible, attributable, and discussable.

Key actions

- Define clear ownership for all AI workloads (models, GPUs, inference endpoints)
- Tag GPU, model, and token usage consistently across environments
- Separate AI spend from general compute costs
- Identify unowned or idle resources

Primary owner

- Engineering + Finance (shared)

Success signals

- 80-90% cost attribution accuracy
- Previously untracked AI spend exposed
- Cost conversations shift from assumptions to facts

Objective

Replace after-the-fact reporting with forward-looking control.

Key actions

- Implement forecasting based on live usage signals
- Incorporate seasonality, release cycles, and workload behaviour
- Enable anomaly detection for unexpected spend
- Route alerts to operational channels, not just finance inboxes

Primary owner

- Finance (with engineering input)

Success signals

- ~90% forecast accuracy
- <5% monthly budget variance
- Cost spikes flagged before invoices arrive

Objective

Reduce waste continuously without manual intervention.

Key actions

- Enable automated rightsizing for underutilised resources
- Enforce schedules for non-production environments
- Test eligible GPU workloads on spot or preemptible capacity
- Eliminate idle or "zombie" resources

Primary owner

- Engineering (within finance-approved guardrails)

Success signals

- 25-35% cumulative savings identified or realised
- Reduced engineering time spent on cost tuning
- No performance or reliability regressions

Objective

Replace after-the-fact reporting with forward-looking control.

Key actions

- Implement forecasting based on live usage signals
- Incorporate seasonality, release cycles, and workload behaviour
- Enable anomaly detection for unexpected spend
- Route alerts to operational channels, not just finance inboxes

Primary owner

- Finance (with engineering input)

Success signals

- ~90% forecast accuracy
- <5% monthly budget variance
- Cost spikes flagged before invoices arrive

Key takeaway

Focus on **speed** over **perfection** in the first 60 days.

Early wins build confidence, alignment, and momentum, unlocking deeper optimisation in later phases.

AI workload cost optimisation cheatsheet

Where AI FinOps savings actually come from

Model selection

Match intelligence to economic value

Not every workload requires frontier intelligence. High-performing teams align model choice to task complexity.

Do this

- Use smaller or specialised models for classification, extraction, routing, and summarisation
- Reserve large models for complex reasoning or customer-facing outputs
- Route requests dynamically based on task complexity

Avoid this

- Defaulting all workloads to the most capable model
- Paying for reasoning depth where it does not change outcomes

Typical impact

- 60–80% cost variance for identical outputs depending on model choice

Token efficiency

Control the invisible multiplier

Token usage compounds quietly. Small prompt changes can double inference costs without obvious signals.

Do this

- Compress prompts and enforce structured outputs
- Set hard caps on context windows and response length
- Trim context dynamically based on task type

Avoid this

- Overly verbose prompts
- Unbounded responses
- Repeated retries without limits

Typical impact

- 30%+ reduction in token consumption with no quality loss

GPU optimisation

Optimise orchestration, not usage

GPUs are the largest and most volatile cost driver in AI stacks. The difference between disciplined and reactive teams is orchestration.

Do this

- Right-size over-provisioned instances
- Enforce shutdown schedules for non-production workloads
- Use spot or preemptible capacity for eligible training jobs
- Separate training and inference architectures

Avoid this

- Leaving GPUs running outside active usage windows
- Treating training and inference as a single cost pool
- Manual scheduling and cleanup

Typical impact

- 70–90% savings on eligible training workloads
- 60%+ idle reduction outside production hours

Key principle

Optimisation should not slow delivery.

When model choice, token limits, and GPU scheduling are automated and policy-driven:

- Engineers stop firefighting AI bills
- Finance gains predictability without blocking experimentation
- Leadership controls margins without operational drag

This is the difference between cutting costs and governing costs.

AI_FinOps_Action_Pack_2026_TardiTech.pdf

Strategic alignment scoring template

A simple decision lens to prioritise FinOps initiatives for AI workloads

Purpose

This template helps leadership, finance, and engineering prioritise AI FinOps initiatives consistently.

Each initiative is scored across four factors to determine what to execute now, what to sequence next, and what to defer.

Factor	Score 1 (low)	Score 2 (medium)	Score 3 (high)	Score (1-3)
Business impact	Marginal or indirect impact Nice-to-have improvement	Supports a strategic pillar Improves internal efficiency or reliability	Direct impact on revenue protection, growth, or customer experience	
Savings potential	Less than 10% savings Unclear or unvalidated upside	10-20% realistic savings potential	Material savings potential (validated or strongly modelled)	
Ease of implementation	High effort Significant engineering time Tooling or architectural changes required	Moderate coordination Some configuration required	Low effort Quick to implement Minimal delivery disruption	
Safety / reversibility	High risk Potential performance regressions High operational uncertainty	Manageable risk with safeguards	Low risk Safe, reversible, well-understood changes	
Total score				____/12

Score interpretation

- 10-12 → Priority initiative (execute now)
- 7-9 → Sequence next (after prerequisites)
- 4-6 → Defer, narrow, or revisit later
- Below 4 → Do not prioritise yet

Optimise deliberately, not aggressively.

Sustainable efficiency comes from sequencing decisions, not chasing every possible saving.

Key metrics tracking matrix

What to measure to govern AI spend effectively

Financial metrics (runway and predictability)				
Metric	Target	Owner	Review cadence	Why it matters
Budget variance	<5% month over month	Finance	Monthly	Signals forecast accuracy and planning confidence
Wasted spend percentage	<10%	Finance + Engineering	Monthly	Indicates maturity of tagging, ownership, and automated cleanup
Cost per customer / feature	Trending down	Leadership	Quarterly	Connects AI and cloud spend to business outcomes
Commitment coverage	70% on steady workloads	Finance	Quarterly	Balances savings with flexibility as usage evolves
Operational metrics (speed and control)				
Metric	Target	Owner	Review cadence	Why it matters
Forecast accuracy	~90%	Finance	Monthly	Confirms predictive controls reflect real usage patterns
Anomaly detection time	<10 minutes	Engineering	Continuous	Prevents small issues from becoming large invoices
Optimisation adoption rate	Trending up	Engineering + Finance	Monthly	Signals trust in FinOps recommendations
Engineering time on cost reviews	Trending down	Leadership	Quarterly	Confirms automation is reducing operational friction
AI-specific metrics (where costs actually hide)				
Metric	Target	Owner	Review cadence	Why it matters
GPU utilisation rate	60-70% sustained	Engineering	Weekly	Highlights idle capacity and scheduling inefficiencies
Token cost per output	Trending down	Engineering	Monthly	Reveals prompt design and model routing inefficiencies
Cost per inference / AI feature	Stable or declining	Product	Monthly	Enables true unit economics for AI-driven features
Training vs inference spend	Controlled	Engineering	Monthly	Prevents experimentation from silently dominating budgets

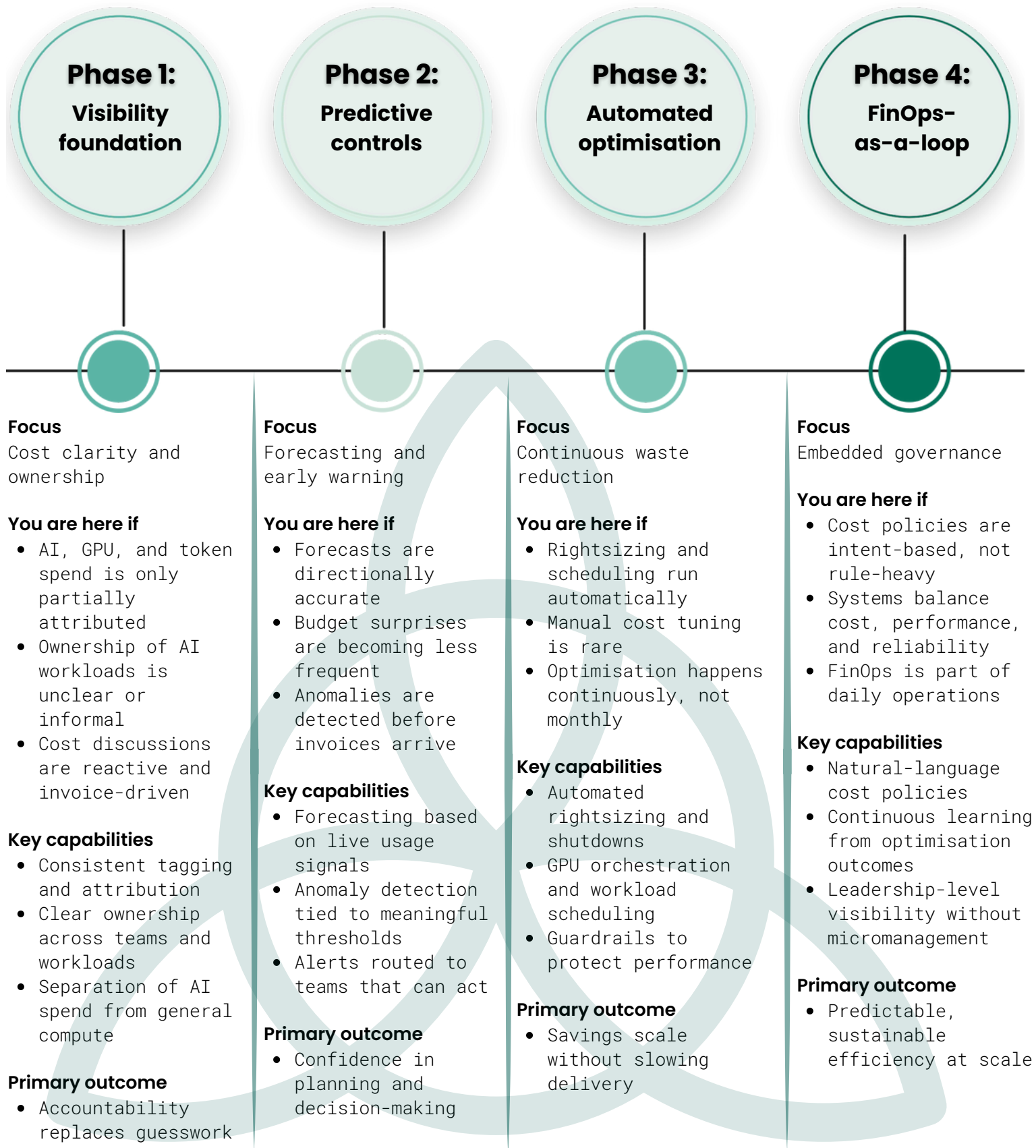
Track fewer metrics. Review them more often. Assign clear ownership.

If a metric does not influence a decision, it does not belong here.

FinOps maturity snapshot

A one-page view of where you are today and what “good” looks like next

This snapshot helps teams understand their current FinOps maturity and what to focus on next.



Key takeaway

The goal is **not** maximum maturity.
The goal is the **right maturity** for your current scale.

AI FinOps readiness checklist

A quick diagnostic before scaling AI further

Confirm foundational controls are in place before expanding AI workloads.

Ownership & accountability

- ☐ Every AI workload (models, GPUs, inference endpoints) has a clear owner
- ☐ Ownership is shared appropriately between engineering and finance
- ☐ Cost accountability is explicit, not assumed

Visibility & attribution

- ☐ AI, GPU, and token spend is tagged consistently
- ☐ AI spend is separated from general compute costs
- ☐ At least 80% of AI-related spend is attributable

Forecasting & predictability

- ☐ AI spend forecasts are directionally reliable
- ☐ Budget variance is consistently under 5%
- ☐ Anomalies are detected before invoices arrive

Automation readiness

- ☐ Rightsizing and shutdowns can be automated safely
- ☐ Non-production workloads follow enforced schedules
- ☐ Guardrails exist to protect performance during optimisation

Governance & decision-making

- ☐ Cost policies are defined in terms of intent, not manual rules
- ☐ Leadership reviews scenarios, not raw spend reports
- ☐ FinOps decisions do not slow engineering delivery

Results interpretation

Mostly checked → Ready to scale AI usage deliberately

Several gaps → Address fundamentals before increasing complexity

Many unchecked → Focus on visibility and ownership first

Key takeaway

Scaling AI without readiness turns cost into risk.
Scaling with readiness turns cost into leverage.

Turn insight into execution

You now have the framework.

The next step is applying it to your environment.

Book a free AI FinOps strategy call

A focused conversation to:

- Assess your current FinOps maturity
- Identify immediate cost and risk exposure
- Prioritise the next 60 days with confidence

No sales pitch. Just practical guidance.

Connect with us

