



Folium

PUBLIC-FACING PDF

REVIEW BEFORE PRODUCTION

FOLIUM SYSTEMS

LOCAL PRIVATE HYBRID AI GUIDE

# Local, Private, And Hybrid AI Guide

Not every AI workload belongs in the same runtime. Folium helps buyers choose cloud APIs, private endpoints, local models, hybrid routes, RAG, agents, and fallback paths based on privacy, cost, latency, portability, control, and support.

## AUDIENCE

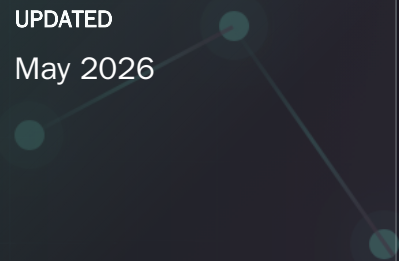
Business owners, technical buyers, security reviewers, local-AI teams, and cost-sensitive operators

## PURPOSE

Explain runtime placement and data control in plain business language

## UPDATED

May 2026



Cloud AI, local AI, private endpoints, and hybrid routing each have a place.

Runtime choice should follow data sensitivity, cost, latency, accuracy, support, and owner review.

Folium can help operate model-agnostic paths across APIs, local models, Ollama, vLLM, SGLang, RAG, agents, databases, and legacy systems.

NEW

BOUNDARY

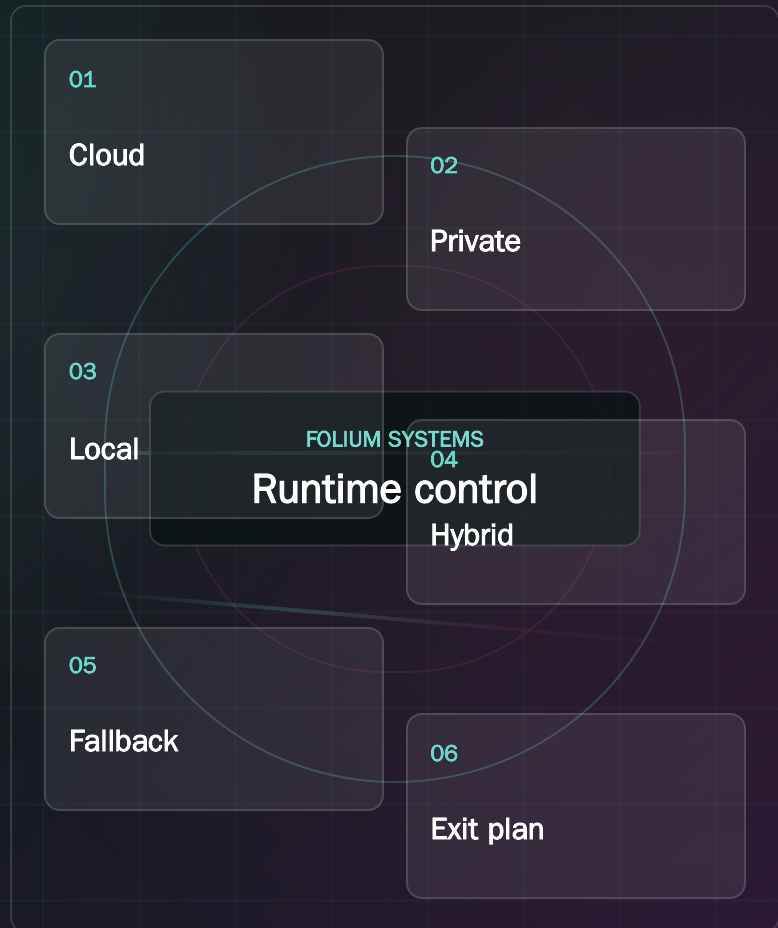
NEXT GATE

## RUNTIME PLACEMENT

# Local, private, and hybrid AI choices should be business decisions.

The runtime guide compares cloud APIs, private endpoints, local models, hybrid routes, RAG placement, agent lanes, costs, latency, custody, and fallback.

## RUNTIME CONTROL



01

Shows Folium can reason beyond one provider.

02

Makes privacy, cost, latency, and control tradeoffs visible.

03

Protects data custody before architecture locks in.

# Choose the review route before reading cover to cover.

This packet is meant to support a real decision meeting. Different reviewers should enter through different routes, then come back together around the same controlled next step.

## DECISION ROUTE

### EXECUTIVE ROUTE

#### Decision first

Start with the cover, visual summary, executive read, controls, first ninety days, and handoff. This route helps leaders decide whether the next move is education, audit, first build, pilot, or operations.

- Outcome
- Risk
- Owner
- Next gate

## OPERATING ROUTE

### OPERATIONS ROUTE

#### How the work will run

Read the workflow map, procedures, operating roles, metrics, first sprint, and buyer worksheet. This route shows whether staff can actually use, review, and improve the future process.

- Workflow
- Staff
- Support
- Improve

## TRUST ROUTE

### TECHNICAL AND TRUST ROUTE

#### Where the boundaries live

Focus on records and work products, controls, risk assumptions, reference work products, source truth, runtime placement, and launch conditions before any private access expands.

- Source
- Access
- Runtime
- Rollback

### BUYER SESSION ROUTE

#### Turn reading into a working session

Use the discovery questions, role review route, buyer worksheet, and engagement fit ladder to prepare one process, one owner, one source map, and one next decision.

- Process
- Examples
- Questions
- Decision

**Best use:** bring one workflow, the people who own it, the systems it touches, the data classes involved, and the decision this packet should help leadership make.

# Local private hybrid AI guide in plain language.

Not every AI workload belongs in the same runtime. Folium helps buyers choose cloud APIs, private endpoints, local models, hybrid routes, RAG, agents, and fallback paths based on privacy, cost, latency, portability, control, and support.

RECORD

BOUNDARY

ACTION

CLOUD

## Cloud APIs for speed and capability

Cloud models can be excellent when data, policy, cost, and latency fit the use case.

- Capability
- Speed
- Vendor

PRIVATE

## Private endpoints for control

Private or dedicated routes can reduce exposure and fit stricter operating rules.

- Boundary
- Access
- SLA

LOCAL

## Local models for custody

Local AI can keep work closer to the business when privacy, cost, or availability matters.

- Ollama
- vLLM
- SGLang

HYBRID

## Hybrid routing for fit

The best architecture may route each task to the runtime that matches its risk and value.

- Router
- Fallback
- Governance

This packet is public-facing. It is written for serious review without exposing private infrastructure, customer data, credentials, live provider wiring, or internal project labels.

# The operating path should be visible before anyone trusts the outcome.

Folium uses workflow maps to turn broad AI ambition into inspectable work. Each phase names the procedure, the visible output, and the decision gate that prevents excitement from outrunning control.

DECISION GRID

REVIEW LENS

NEXT STEP

| PHASE                 | PROCEDURE  | VISIBLE OUTPUT              | DECISION GATE                               |
|-----------------------|--|-----------------------------|---|
| <b>Classify data</b>  | Identify public, internal, confidential, regulated, secret, and blocked data.          | Data sensitivity map.       | Runtime cannot be chosen before data class. |
| <b>Classify tasks</b> | Separate drafting, retrieval, reasoning, summarization, routing, and action tasks.     | Task authority map.         | Each job has a role.                        |
| <b>Choose runtime</b> | Compare cloud, private, local, hybrid, and manual fallback for each task.              | Runtime placement matrix.   | The route fits the job.                     |
| <b>Design RAG</b>     | Define source ingestion, chunking, metadata, retrieval, freshness, and evaluation.     | Knowledge architecture.     | Answers stay source-grounded.               |
| <b>Design agents</b>  | Set tool permissions, human approval, logs, blocked actions, and escalation.           | Agent control spec.         | Agents cannot exceed authority.             |
| <b>Deploy safely</b>  | Prepare containers, services, monitoring, API boundaries, cost controls, and rollback. | Deployment guide.           | Operations can support it.                  |
| <b>Evaluate</b>       | Test accuracy, retrieval, latency, cost, privacy, failure handling, and fallback.      | Evaluation report.          | Runtime earns trust.                        |
| <b>Operate</b>        | Monitor usage, update sources, rotate models, review cost, and maintain support.       | AI runtime operations loop. | The route stays useful.                     |

# The work should leave behind material a buyer can inspect.

A serious engagement should produce more than conversation. Folium packages records, diagrams, checklists, routes, system surfaces, launch gates, and handoff material so the buyer can keep control after the first win.

DECISION GRID

REVIEW LENS

NEXT STEP

| WORK PRODUCT                    | WHAT IT CONTAINS   | HOW THE REVIEWER USES IT         |
|---------------------------------|--|----------------------------------|
| <b>Runtime placement matrix</b> | Tasks mapped to cloud, private, local, hybrid, or manual fallback.               | Prevents one-size-fits-all AI.   |
| <b>Data custody map</b>         | Where data originates, moves, persists, and gets blocked.                        | Supports security review.        |
| <b>Model route catalog</b>      | Candidate providers, local models, endpoints, cost, latency, and quality notes.  | Keeps the system model-agnostic. |
| <b>RAG architecture</b>         | Sources, embeddings, retrieval, metadata, freshness, and evaluation.             | Protects answer quality.         |
| <b>Agent permission spec</b>    | Tools, allowed actions, approvals, blocked actions, logs, and escalation.        | Controls automation risk.        |
| <b>Operations guide</b>         | Monitoring, support, rollback, incident path, model updates, and source refresh. | Makes runtime supportable.       |

# The procedure is the product as much as the technology.

The goal is not to make AI look impressive for one meeting. The goal is to make the operating path repeatable, explainable, reviewable, and safe enough to improve.

## CHECKLIST

## OWNER PATH

## RELEASE SIGNAL

- Do not choose local or cloud by ideology; choose by task, data, cost, and support need.
- Classify data before moving it into any model route.
- Separate retrieval from generation and action.
- Keep sensitive workflows close to approved environments.
- Use local models where custody, cost, offline access, or control matters.
- Use cloud APIs where capability, speed, or managed service fit the risk.
- Use hybrid routing when different tasks need different controls.
- Measure latency, cost, quality, support burden, and fallback behavior.
- Define who can approve model changes.
- Keep a model and route lifecycle record.

# Governance, quality, and launch gates keep speed honest.

Folium keeps the buyer's next decision tied to observable gates: source truth, authority, access, testing, ownership, support, rollback, and improvement cadence.

DECISION GRID

REVIEW LENS

NEXT STEP

| GATE                | WHAT MUST BE TRUE                                       | STOP OR REFINE SIGNAL                    |
|---------------------|---|--|
| <b>Data gate</b>    | Sensitivity, retention, and custody rules are known.    | Data class is unknown.                   |
| <b>Runtime gate</b> | Each task has a justified route and fallback.           | Everything goes to one model by default. |
| <b>Cost gate</b>    | Expected usage, token/compute cost, and alerts exist.   | Spend can grow silently.                 |
| <b>Quality gate</b> | Retrieval, generation, and task success are evaluated.  | The route is chosen only by preference.  |
| <b>Support gate</b> | Monitoring, updates, incidents, and rollback are owned. | No one can operate the runtime.          |

# The right questions expose the real project.

These prompts help a buyer and Folium decide whether the next step should be education, audit, first build, security review, pilot, or an operating support path.

CHECKLIST

OWNER PATH

RELEASE SIGNAL

- What data cannot leave your environment?
- What tasks require the strongest model and which require a controlled local route?
- What latency is acceptable?
- What cost would become painful at scale?
- Who owns model updates and source refresh?
- Which tasks need RAG instead of general model memory?
- Which agent actions require human approval?
- What happens when the preferred model route fails?

# Diagrams, charts, and overlays make the work easier to review.

Dense AI work should not only be explained in paragraphs. The reviewer should be able to inspect maps, scorecards, matrices, lanes, and before-after views that reveal where the value and risk live.

RECORD

BOUNDARY

ACTION

## Runtime router

A routing diagram from task class and data class to cloud, private, local, hybrid, or manual fallback.

- Task
- Data
- Route
- Fallback

## RAG pipeline

A pipeline from source intake to chunking, metadata, embedding, retrieval, answer, and evaluation.

- Source
- Index
- Retrieve
- Evaluate

## Cost-control chart

A chart comparing API token spend, local compute, support burden, and expected usage.

- Tokens
- Compute
- Support
- Scale

## Model lifecycle board

A board for active, experimental, parked, retired, and replacement routes.

- Active
- Test
- Park
- Retire

# Every serious AI path needs named owners before it becomes dependency.

The same technology can be safe or unsafe depending on who owns the workflow, data, quality, launch authority, support, and improvement loop. Folium makes those responsibilities explicit so no buyer inherits an orphaned system.

DECISION GRID

REVIEW LENS

NEXT STEP

| ROLE                             | OWNS   | RECORD TO INSPECT   |
|----------------------------------|--|---|
| <b>Executive sponsor</b>         | Priority, budget, risk tolerance, stop/continue decision, and expansion timing.          | Decision note, value hypothesis, and approval boundary.     |
| <b>Business process owner</b>    | The day-to-day work, acceptance criteria, staff impact, and operational usefulness.      | Workflow map, user feedback, and adoption notes.            |
| <b>Technical owner</b>           | Systems, APIs, databases, runtime placement, deployment, monitoring, and fallback.       | Architecture map, integration log, and support route.       |
| <b>Knowledge owner</b>           | Source truth, document freshness, policies, retrieval scope, and correction workflow.    | Source inventory, freshness cadence, and review exceptions. |
| <b>Security or risk reviewer</b> | Data classes, credentials, access, logs, retention, blocked actions, and incident path.  | Boundary map, permission table, and rollback trigger.       |
| <b>Folium delivery lead</b>      | Build coordination, review file, known limits, quality checks, and handoff completeness. | Launch room, eval record, and improvement backlog.          |

# A max-detail packet should tell reviewers how to judge the work.

Folium uses scorecards to make a subjective AI conversation more inspectable. The score is not a substitute for judgment; it helps leadership see whether the next step is education, repair, sandbox, pilot, or operations.

DECISION GRID

REVIEW LENS

NEXT STEP

| SCORE AREA                  | STRONG SIGNAL  | WEAK SIGNAL   |
|-----------------------------|--|---|
| <b>Business fit</b>         | The workflow is specific, painful, owned, and tied to measurable operational improvement.      | The project is framed as adding AI generally.           |
| <b>Source truth</b>         | Approved sources are known, fresh, classified, and connected to the answer path.               | The system mixes stale, unknown, or unapproved sources. |
| <b>Behavior quality</b>     | Representative tasks pass, wrong-answer behavior is known, and edge cases are recorded.        | The review build only shows a polished happy path.      |
| <b>Authority control</b>    | AI actions are separated into draft, retrieve, recommend, route, execute, block, and escalate. | The system can act without visible permission.          |
| <b>Staff readiness</b>      | Users can explain the tool, correct it, escalate, and understand their role.                   | Staff feel replaced, confused, or unsupported.          |
| <b>Operations readiness</b> | Support, monitoring, rollback, release rhythm, and source refresh are owned.                   | No one knows who maintains the system after launch.     |

# The work should have a believable first ninety days.

A controlled first ninety days keeps ambition high without turning uncertainty into production risk. Folium uses the period to move from understanding into a narrow working example, then into reviewable operating rhythm.

DECISION GRID

REVIEW LENS

NEXT STEP

| WINDOW               | FOCUS   | EXPECTED OUTPUT  |
|----------------------|---|--|
| <b>First 30 days</b> | Discovery, source inventory, first-lane selection, staff interviews, data boundary, and build plan.             | Process map, owner map, first-build scope, source list, and launch blockers.   |
| <b>Days 31-60</b>    | Working surface, RAG or agent behavior, integration stub, evaluation cases, browser checks, and staff review.   | Sandbox, evaluation file, screenshots, known limits, and repair list.          |
| <b>Days 61-90</b>    | Architecture review, pilot conditions, governance layer, training guide, support path, and improvement cadence. | Launch room, go/no-go record, operations guide, and next-stage recommendation. |

# The hidden assumptions should be visible before they become expensive.

Every AI engagement contains assumptions about data, people, systems, cost, behavior, and authority. Folium treats those assumptions as review material, not background noise.

DECISION GRID

REVIEW LENS

NEXT STEP

| ASSUMPTION                         | WHY IT MATTERS  | HOW FOLIUM REVIEWS IT   |
|------------------------------------|---|---|
| <b>The source is authoritative</b> | AI can only be as reliable as the sources and business rules it is allowed to use.                      | Source inventory, owner confirmation, retrieval tests, freshness cadence.   |
| <b>The process is ready</b>        | A broken process can become a faster broken process when AI is added too early.                         | Workflow mapping, bottleneck review, owner interview, first-lane narrowing. |
| <b>The runtime fits the data</b>   | Cloud, private, local, and hybrid routes carry different privacy, cost, latency, and support tradeoffs. | Runtime matrix, data classification, provider review, fallback plan.        |
| <b>Staff will adopt the tool</b>   | Adoption fails when users do not understand, trust, correct, or benefit from the system.                | Training notes, staff review, feedback loop, manager visibility.            |
| <b>Authority is clear</b>          | The system can create harm if it sends, updates, approves, or routes without permission.                | Permission table, blocked actions, human review, audit trail.               |
| <b>The system can be supported</b> | A useful first build becomes fragile if nobody owns incidents, source updates, or cost review.          | Support guide, owner map, release rhythm, rollback trigger.                 |

# The first sprint should produce something real and reviewable.

Folium prefers a narrow first sprint that creates a working surface or review file the buyer can challenge. The first sprint is not the final system; it is the safest way to make the future visible.

## CHECKLIST

## OWNER PATH

## RELEASE SIGNAL

- Confirm the single process and the decision the sprint must support.
- Collect approved example material, redacted review records, public references, screenshots, workflow notes, and source rules.
- Define what will be built: portal, dashboard, RAG assistant, agent route, integration adapter, audit file, or launch room.
- Create the visual workflow: intake, source, model or agent route, human review, output, record, and next gate.
- Run representative tasks, edge cases, bad input, missing data, and blocked-action tests.
- Prepare browser screenshots, known limits, support questions, and next-stage blockers.
- Review with staff and leadership before expanding data, access, authority, or dependency.
- End with a decision: stop, refine, rebuild, pilot, or prepare an operating plan.

# The packet should make the invisible work tangible.

AI work often fails because the important pieces are invisible until something breaks. Folium turns those pieces into work products the buyer can open, print, challenge, and improve.

## RECORD

## BOUNDARY

## ACTION

## Process map

A before-and-after workflow showing people, systems, data, decision points, blockers, and expected output.

- Before
- After
- Owner
- Gate

## Data boundary map

A map of source classes, approved use, blocked use, retention, provider exposure, and custody.

- Public
- Internal
- Private
- Blocked

## Model and agent route

A path showing which model, tool, retrieval source, or agent lane is used and where humans approve.

- Route
- Tool
- Review
- Escalate

## Evaluation file

A record of tasks, expected outcomes, failures, repairs, known limits, and acceptance criteria.

- Cases
- Failures
- Repairs
- Limits

## Launch room

A board for owners, support, training, rollback, incidents, go/no-go, and improvement backlog.

- Owner
- Support
- Rollback
- Backlog

## Handoff guide

A plain-language guide staff can use to understand what the system does, cannot do, and how to report problems.

- Use
- Limit
- Correct
- Report

# The business should know how improvement will be measured.

Folium keeps measurement practical. The first goal is not a perfect dashboard; it is a clear set of signals that shows whether the process is saving time, reducing risk, strengthening staff, or improving customer outcomes.

DECISION GRID

REVIEW LENS

NEXT STEP

| SIGNAL                  | WHAT TO WATCH   | DECISION IT SUPPORTS   |
|-------------------------|---|--|
| <b>Time recovered</b>   | Manual steps removed, average handling time, repeated work reduced, faster routing.           | Should this workflow expand to more users or adjacent processes? |
| <b>Quality improved</b> | Wrong answers, missing sources, correction rate, review exceptions, customer rework.          | Is behavior strong enough for pilot or does it need repair?      |
| <b>Risk reduced</b>     | Blocked unsafe actions, escalations, data-boundary violations avoided, rollback readiness.    | Can authority expand or should controls remain tight?            |
| <b>Staff confidence</b> | Training completion, feedback volume, adoption friction, override rate, manager notes.        | Does the workforce need more support before launch?              |
| <b>Cost and runtime</b> | Provider cost, local infrastructure cost, latency, uptime, fallback use, subscription sprawl. | Should runtime placement change?                                 |
| <b>Customer impact</b>  | Response speed, consistency, issue resolution, conversion support, satisfaction signals.      | Is the capability improving the business outcome?                |

# Each reviewer should know what to inspect first.

A max-detail packet is only useful when different reviewers can find their lane quickly. Folium separates executive, operations, technical, security, finance, and staff questions so the buyer can bring the right people into the right part of the review.

DECISION GRID

REVIEW LENS

NEXT STEP

| REVIEWER                         | START WITH   | DECISION THEY SUPPORT   |
|----------------------------------|--|---|
| <b>Executive sponsor</b>         | Value hypothesis, launch gate, first ninety days, and stop/refine/continue choices.      | Whether the process deserves a controlled engagement.         |
| <b>Operations lead</b>           | Workflow map, operating roles, support rhythm, and staff feedback loop.                  | Whether the future process can be run by the team.            |
| <b>Technical lead</b>            | Runtime placement, data path, integration surface, monitoring, and fallback.             | Whether the architecture can be supported safely.             |
| <b>Security or risk reviewer</b> | Data classes, permissions, blocked actions, logs, retention, and rollback.               | Whether access can expand beyond public review.               |
| <b>Finance or owner</b>          | Cost signals, subscription overlap, runtime tradeoffs, labor impact, and support burden. | Whether the first build has a practical business case.        |
| <b>Staff user</b>                | Plain-language use, limits, escalation, correction path, and training expectations.      | Whether the tool strengthens the job instead of confusing it. |

# The packet should turn into a working session, not only reading material.

Before a call, Folium wants the buyer to gather the real operating pieces that make the review useful. The worksheet keeps the conversation grounded in one process, one owner, one source map, and one next decision.

## CHECKLIST

## OWNER PATH

## RELEASE SIGNAL

- Bring one workflow that is slow, risky, expensive, repetitive, customer-visible, or staff-heavy.
- Name the systems touched by the workflow: store, CRM, ERP, inbox, spreadsheet, database, portal, document folder, or legacy application.
- Separate approved public material from internal, customer, regulated, confidential, credential, and blocked material.
- Write down who owns the work today, who reviews exceptions, and who will own the AI-assisted version.
- List the decisions AI may draft, retrieve, recommend, route, block, or escalate, and the decisions that stay human-owned.
- Bring examples of good output, bad output, common exceptions, missing data, and customer-facing risk.
- Name the first useful working surface: dashboard, portal, assistant, queue, control room, commerce lane, integration, or review file.
- Decide what record would make leadership comfortable with the next stage.

# The next step should match the maturity of the record.

Folium does not need every buyer to start at the same altitude. The right offer depends on how much process clarity, source truth, owner alignment, and launch readiness already exists.

DECISION GRID

REVIEW LENS

NEXT STEP

| IF THE BUYER HAS  | BEST NEXT FOLIUM MOVE                          | OUTPUT TO EXPECT  |
|---|--|---|
| <b>AI interest but no clear process</b>                 | AI systems audit or first workflow finder.     | Pressure map, source inventory, first-lane recommendation, and risk view.                 |
| <b>A clear process but no working surface</b>           | Forward engineering first sprint.              | Clickable surface, route map, known limits, and next-stage blockers.                      |
| <b>A tool that works in parts but not in operations</b> | Architecture and launch readiness review.      | Permission map, runtime decision, support model, and go/no-go record.                     |
| <b>A failed or frightening rollout</b>                  | AI recovery and staff enablement path.         | Issue register, staff training plan, repair roadmap, and confidence loop.                 |
| <b>Sensitive data or cost pressure</b>                  | Local, private, or hybrid AI placement review. | Runtime matrix, data custody plan, fallback route, and vendor-exit view.                  |
| <b>A useful pilot that needs care</b>                   | AI operations support.                         | Monitoring rhythm, source refresh, release notes, incident path, and improvement backlog. |

# The last page of a packet should create the next controlled move.

Folium's handoff view separates what can be done now, what needs customer records, what needs approval, and what should wait until the review file is stronger.

DECISION GRID

REVIEW LENS

NEXT STEP

| HANDOFF LANE     | OWNER                         | NEXT RECORD   |
|------------------|-------------------------------|---|
| Data owner       | Security or operations        | Data custody and sensitivity map.                       |
| Runtime owner    | Technical lead                | Runtime placement and deployment guide.                 |
| Quality owner    | Folium and customer reviewers | Evaluation results and known limits.                    |
| Operations owner | Support lead                  | Monitoring, source refresh, and model lifecycle record. |

The strongest next step is narrow: one process, one owner, one source map, one working surface, one review file, and one decision gate.

# Private AI is not a slogan. It is a placement decision with operating consequences.

Use this guide to decide which workloads should use cloud, private, local, hybrid, RAG, agent, or manual fallback paths.

## Bring the process

Name the business process, the systems involved, the people affected, and the decision this PDF should support.

## Separate review from production

Keep public examples, sandbox review, pilot access, and production dependency in separate stages with clear owners.

## Ask for the record

Request screenshots, browser checks, known limits, launch blockers, support plans, and the next approval path.