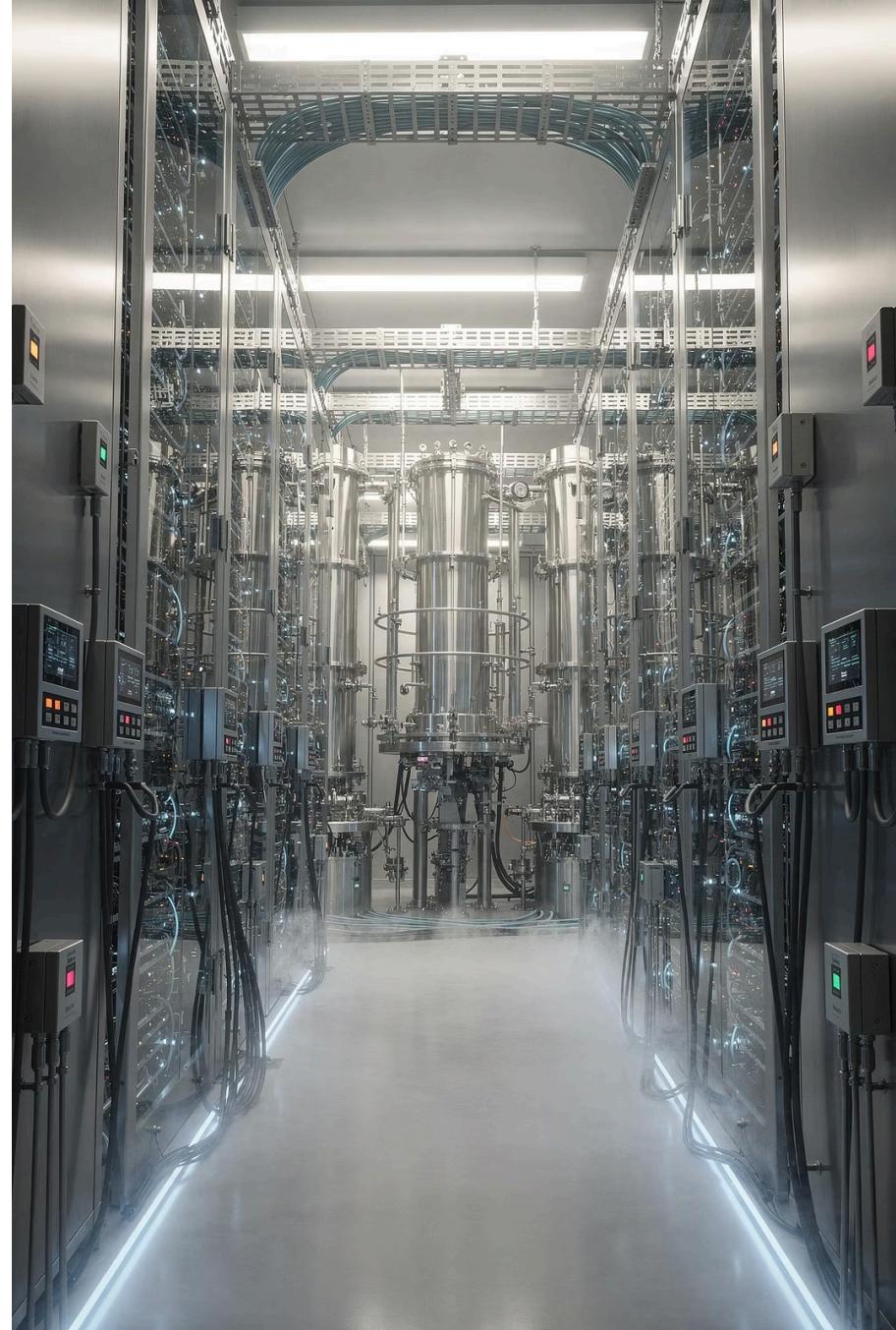


Quantum Sidecar Architecture for Plant Sensor Root-Cause Analysis

A comprehensive guide to building quantum-adjacent Python services for industrial anomaly diagnosis using Cursor, Qiskit, and IBM Quantum hardware. This technical blueprint demonstrates how AI-assisted development accelerates quantum application delivery whilst maintaining architectural rigour and engineering best practices.



Why Cursor Excels for Quantum Development Work



AI-First Repo-Wide Edits

Cursor is an AI-powered IDE built atop VS Code that comprehends and refactors across your entire plant-sensor-quantum repository. This capability proves invaluable when iterating simultaneously on QUBO models, Qiskit quantum circuits, and FastAPI service layers. The AI assistant understands contextual relationships between modules, enabling sophisticated cross-file modifications that would require substantial manual coordination in traditional editors.



Quantum-Adjacent Python Excellence

Most quantum SDKs—including Qiskit, PennyLane, and Classiq client libraries—are fundamentally Python-centric. Cursor handles Python ergonomically with intelligent completions, type-aware refactoring, and seamless integration with quantum frameworks. The editor has demonstrated success with CUDA-Q-style stacks and other quantum computing toolchains, making it a proven choice for hybrid classical-quantum development workflows.



Vibe-Coding Workflow Integration

Cursor's code-aware chat and "apply diff" workflow map directly to vibe-coding use cases: translating natural-language descriptions of plant anomalies or QUBO formulations into concrete implementations, then presenting diffs for safe review. This approach bridges the gap between conceptual quantum algorithms and production-ready code, enabling rapid prototyping without sacrificing code quality or architectural coherence.

Core Technology Stack Within Cursor

Primary Development Components

Qiskit forms the foundation for QAOA implementation, QUBO-to-Ising encoding, and IBM Runtime integration. As the de-facto standard for gate-based quantum programming on IBM hardware, Qiskit provides rich APIs for variational algorithms, optimization primitives, and quantum circuit construction.

PennyLane offers optional support for quantum machine learning experiments. Its differentiable programming model and hardware-agnostic interface enable exploration of hybrid quantum-classical models when QML approaches might complement QAOA-based root-cause analysis.

FastAPI powers the sidecar service API layer. This modern Python web framework delivers high performance, automatic OpenAPI documentation, and native async support—critical for managing quantum job submissions and result retrieval from IBM backends. The framework's dependency injection system elegantly separates concerns between API endpoints, business logic, and quantum execution layers.

pytest provides comprehensive test coverage across unit tests (pure QUBO logic), integration tests (end-to-end API flows), and quantum simulator validation. The framework's fixture system naturally models quantum backend mocking and test data generation for anomaly scenarios.



- ❑ **Development Environment:** Python 3.10+ with virtual environment isolation ensures reproducible builds. Type hints and mypy static analysis enforce interface contracts between classical and quantum code boundaries.

Day-to-Day Development Workflow

01

Repository Initialization

Begin from your defined repository skeleton and open the project in Cursor. Ensure your virtual environment activates automatically and dependencies install cleanly. Configure Cursor's workspace settings to recognise quantum-specific imports and enable AI features for the entire repository scope.

02

Interactive Development with Cursor Chat

Leverage Cursor's AI chat to scaffold implementations from architectural specifications. For example, prompt "Flesh out psq/qubo/model.py from the FMEA topology documentation" to generate initial QUBO construction logic, or "Generate psq/quantum/qaoa_solver.py using Qiskit's QAOA primitives with proper error handling and logging" to create the quantum solver module with production-ready patterns.

03

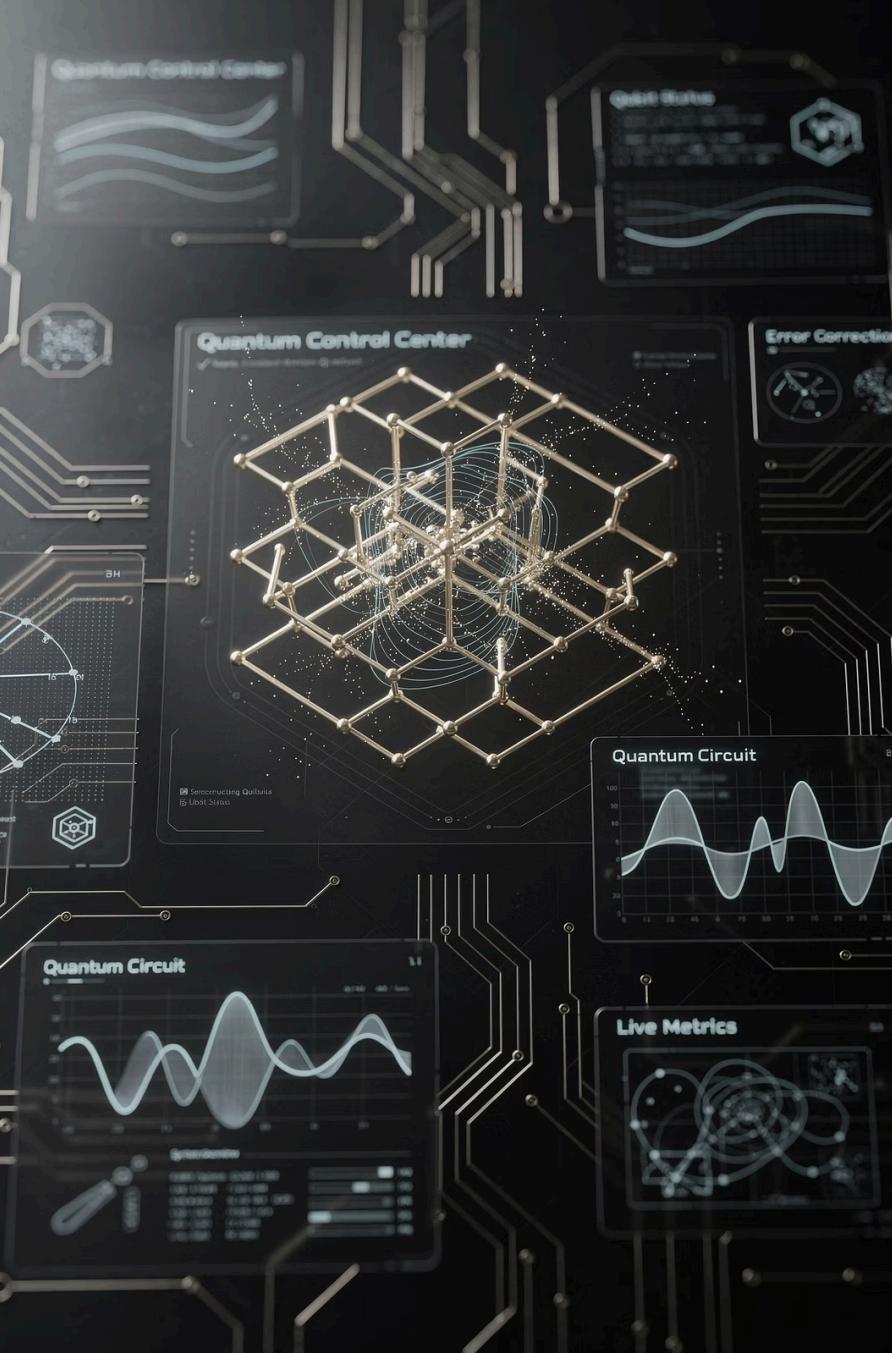
Service Layer Development

Scaffold fastapi_app.py and test clients through conversational prompts that reference your API contracts. The AI assistant generates FastAPI route handlers, Pydantic request/response models, and pytest fixtures for endpoint testing. Request "Create POST /diagnose-plant-anomaly with full request validation and structured logging" to build production-ready service endpoints.

04

Iterative Refinement

Execute unit tests and Jupyter notebooks locally to validate quantum behaviour. When requirements evolve—such as adding new pattern constraint types or modifying the QUBO energy function—describe changes in natural language. Cursor generates appropriate diffs across psq/qubo/, psq/quantum/, and test modules whilst preserving architectural boundaries you've established.



Complementary Tooling Ecosystem

Whilst Cursor serves as your primary development environment, **IBM's Qiskit Code Assistant** provides specialised quantum code completions and templates. This VS Code extension integrates seamlessly with Cursor since it inherits the broader VS Code ecosystem. The assistant offers context-aware suggestions for quantum circuit construction, operator definitions, and IBM Runtime API patterns.

Consider installing the Qiskit Code Assistant for enhanced quantum-specific features including circuit visualisation, quantum gate suggestion, and runtime primitive templates. The combination of Cursor's general AI capabilities with Qiskit's domain-specific intelligence creates a powerful development environment optimised for quantum-adjacent services.

Cursor: AI-Assisted Coding with Architectural Control

The statement "Cursor gives you AI-assisted coding directly inside a full IDE whilst you stay in control of architecture and reviews" captures the essence of modern quantum development workflows. This section unpacks why this combination delivers exceptional results for quantum sidecar services.

Why Python + Qiskit + Cursor Works Brilliantly

Unified Stack Benefits

Qiskit serves as the de-facto standard for gate-based quantum programming on IBM hardware, offering comprehensive APIs for QAOA, VQE, and QUBO/Ising workflows. Python's ubiquity across quantum SDKs and enterprise integration layers (SAP, HPC orchestration, microservices) means one technology stack addresses both quantum algorithms and classical orchestration concerns seamlessly.

Human-in-the-Loop AI

Cursor's conversational interface and diff-based approval model enable rapid generation of Qiskit code, module refactoring, and FastAPI endpoint wiring from natural-language specifications. Crucially, developers approve every change, maintaining code review rigour. This workflow accelerates development velocity without sacrificing quality gates or architectural governance.

Quantum Sidecar Fit

Quantum algorithm code—QUBO modelling, ansatz design, QAOA parameterisation—evolves rapidly during research and optimisation phases. AI assistance accelerates these iterations substantially. Simultaneously, integration with SAP systems, structured logging, and performance benchmarking demands robustness; comprehensive tests, type hints, and manual reviews ensure production reliability whilst Cursor handles boilerplate.

Architecture-Driven Development with Cursor

Designing architecture "with the help of Cursor" means employing the AI assistant as an *interactive collaborator* whilst retaining ownership of architectural decisions. You sketch high-level structure, technical constraints, and design patterns; Cursor materialises and refactors implementations to match your specifications. This approach combines the speed of AI-generated code with the rigour of human architectural oversight.

The methodology treats Cursor as an **architectural executor** rather than architect: you control strategic technology choices, module boundaries, and quality standards, delegating mechanical coding tasks and cross-cutting refactoring work to the AI assistant. This division of labour optimises both development velocity and system quality.



Step 1: Establish Written Architecture Specification

Create ARCHITECTURE.md

Author a top-level ARCHITECTURE.md document (or comprehensive design specification) within your repository that serves as the canonical source of architectural truth. This document should explicitly define module responsibilities, key data flows, and technology selections.

For the quantum sidecar example, document the complete pipeline: "Anomaly JSON → QUBO construction → QAOA execution on IBM Quantum → Root-cause JSON response". Specify that Python, Qiskit, and FastAPI form the core stack, and outline your testing philosophy (unit tests for QUBO logic, integration tests for service endpoints, simulator validation for quantum behaviour).

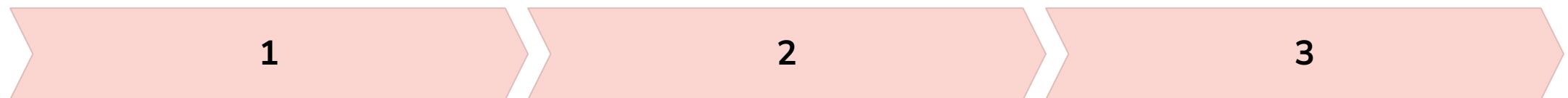
Define Package Structure

Detail the directory layout and module organisation using the psq/ package hierarchy established earlier. Clarify which modules handle data ingestion (psq/data/), QUBO formulation (psq/qubo/), quantum execution (psq/quantum/), orchestration logic (psq/service/), and API exposure (psq/api/).

This specification prevents architectural drift by establishing clear boundaries between concerns. When Cursor generates code, it references this structure to maintain consistency. As architecture evolves, updating ARCHITECTURE.md becomes the single source of truth that drives code changes across the repository.

Step 2: Scaffold Repository Structure with Cursor

With your ARCHITECTURE.md established, leverage Cursor to materialise the skeletal codebase. This step transforms architectural specification into tangible package structure, preparing the foundation for subsequent implementation work.



Create Directory Layout

Establish folders matching your architectural design: `src/psq/` with appropriate subpackages (`data/`, `qubo/`, `quantum/`, `service/`, `api/`), alongside `tests/`, `notebooks/`, and `docs/` directories. Cursor can generate this structure from a simple prompt referencing `ARCHITECTURE.md`.

Generate Module Skeletons

In Cursor's chat interface, prompt: "Given `ARCHITECTURE.md`, create minimal `__init__.py` files and empty module stubs for all packages following the described structure. Add type-annotated function skeletons with docstrings only, no implementation logic yet." This generates interface contracts without premature implementation details.

Review and Refine

Examine generated files, adjusting naming conventions, function signatures, and module boundaries as needed. This review ensures the skeleton matches your architectural vision before implementation begins. Commit this skeleton as your architectural baseline.

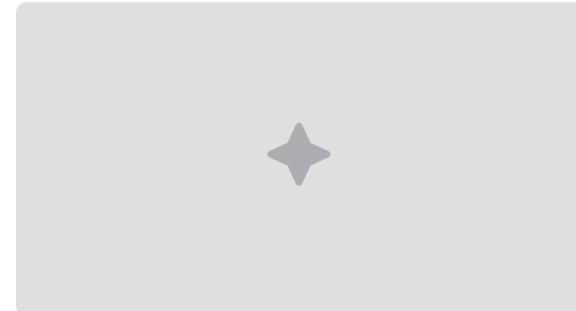
Step 3: Implement Modules Top-Down

Structured Implementation Approach

Work through modules systematically, one responsibility domain at a time, always anchoring prompts to architectural specifications. This discipline prevents scope creep and maintains clean separation of concerns.

QUBO Model Example: "Open psq/qubo/model.py. Implement build_root_cause_qubo as described in ARCHITECTURE.md: accept a list of abnormal sensors and root-cause patterns, produce a QUBO dictionary and variable index mapping. Use the energy function from this LaTeX specification: [insert energy function]. Include comprehensive tests in tests/test_qubo_model.py with toy problem instances."

QAOA Solver Example: "In psq/quantum/qaoa_solver.py, implement run_qaoa_root_cause using Qiskit's QAOA class. Accept a SparsePauliOp cost operator, backend configuration, and QAOA depth parameter. Return bitstring counts with metadata. Follow logging patterns from config.py and logging_utils.py established project-wide."



 **Best Practice:** Reference ARCHITECTURE.md explicitly in every prompt. This grounds Cursor's generated code in your architectural decisions, ensuring consistency across modules developed over time.

Step 4: Enforce Architecture Through Constraints

Prompt-Level Constraints

Embed architectural rules directly in Cursor prompts: "Do not import FastAPI or any web framework dependencies in core psq/ modules—HTTP concerns belong exclusively in psq/api/." Or: "Use type hints for all function signatures; prohibit global mutable state; configuration must flow exclusively through config.py dependency injection."

These explicit constraints guide code generation towards architecturally sound implementations, preventing common anti-patterns like tight coupling between layers or leaky abstractions.

Test-Driven Architecture Validation

Back prompt constraints with comprehensive tests using pytest for behavioural validation. Implement optional static analysis with mypy for type safety and ruff for code quality enforcement. Configure these tools in continuous integration pipelines to catch architectural violations automatically.

When architectural drift occurs—for example, quantum execution code inappropriately appearing in API route handlers—prompt Cursor explicitly: "Refactor so API endpoints only invoke psq/service/orchestrator.py. Move all Qiskit imports into psq/quantum/ and adjust imports across affected modules." Cursor excels at these cross-file architectural refactoring tasks.

Step 5: Evolve Flows with High-Level Prompts

Once your skeletal implementation demonstrates end-to-end functionality on toy examples, use Cursor to introduce new features and optimisations without compromising architectural integrity. High-level prompts that reference architecture documentation enable sophisticated evolution whilst maintaining system coherence.

Feature Addition Example

"Add support for a new QuboRootCauseRequest field called equipment_graph representing physical topology constraints between sensors. Only psq/qubo/model.py and psq/data/schemas.py should require changes to accommodate this field. Update affected unit tests to verify topology constraint handling in QUBO energy functions."

This prompt constrains the change scope explicitly, preventing unintended modifications to unrelated modules and maintaining architectural boundaries.

Performance Optimisation Example

"Optimise run_qaoa_root_cause to support both Aer simulators and IBM Runtime backends interchangeably, controlled by a backend_config object passed at runtime. The public function signature must remain stable for backward compatibility. Add backend selection logic internally with appropriate error handling for unavailable backends."

This approach introduces performance enhancements whilst preserving API contracts and architectural layering. Cursor generates the implementation details whilst you control the architectural evolution.

Architecture-Driven Development Summary

You design the architecture through prose specifications and high-level structural diagrams. Cursor serves as your implementation partner, handling:

- Scaffolding packages and modules conforming to architectural specifications
- Implementing functions under your explicit constraints (technology stack, coding style, module boundaries, testing requirements)
- Executing refactoring operations when architecture evolves, maintaining consistency across the codebase

Employed this way, Cursor becomes an **architectural executor** rather than architect: you retain control over strategic decisions, delegating repetitive mechanical tasks and cross-cutting code transformations to AI assistance.



Using LLMs to Draft ARCHITECTURE.md

An LLM proves exceptionally useful for drafting and iterating on ARCHITECTURE.md, particularly when you've already conceptualised the high-level pattern (plant sensor QUBO → QAOA on IBM → root-cause output) but need to materialise comprehensive documentation. Below we explore an exemplar ARCHITECTURE.md for the plant-sensor-quantum project—precisely the type of document you can co-create and refine collaboratively with Cursor or similar AI assistants.

The following sections present a complete architectural specification that serves as both human-readable documentation and machine-parseable guidance for AI-assisted implementation. This document becomes the single source of truth that drives all subsequent development work.

ARCHITECTURE.md: Goal and Scope



Service Purpose

This quantum sidecar service accepts plant sensor anomaly windows alongside a library of known root-cause patterns, formulates a combinatorial optimisation problem as a QUBO (Quadratic Unconstrained Binary Optimisation), solves it using QAOA (Quantum Approximate Optimisation Algorithm) on IBM Quantum hardware or simulators, and returns a ranked list of candidate root causes with their associated sensor coverage.



Scope Definition

The service focuses exclusively on *combinatorial root-cause assignment*, not initial anomaly detection (which remains a classical preprocessing concern). It provides a clean HTTP API suitable for integration with SAP Asset Intelligence Network (AIN), Plant Maintenance (PM), Manufacturing Execution Systems (MES), or edge computing agents. The architecture supports both IBM Quantum backends and local Qiskit Aer simulators for development and testing workflows.

High-Level Data Flow Architecture

Input Acquisition

Upstream systems (SAP AIN/PM/MES or dedicated anomaly detection services) invoke POST /diagnose-plant-anomaly with structured JSON payload containing: anomaly window identifier and metadata, list of abnormal sensors with severity scores (e.g., z-scores or residual magnitudes), candidate root-cause patterns with their sensor mappings and optional topology constraints.

Quantum Optimisation

The QUBO undergoes conversion to an Ising Hamiltonian, then transformation into a Qiskit SparsePauliOp suitable for quantum execution. QAOA runs via IBM Qiskit Runtime (for real quantum hardware) or local Aer simulators (for development), approximately minimising the cost function. Resulting bitstring samples encode candidate root-cause assignments.

QUBO Construction

The service constructs a QUBO energy function representing trade-offs between: maximising coverage of severe anomalies, minimising the number of selected root-cause patterns (parsimony principle), and maintaining consistency with known pattern–sensor relationships derived from FMEA documentation and equipment topology graphs.

Result Post-Processing

The service decodes bitstring samples into structured JSON responses containing: ranked root-cause hypotheses with confidence scores, coverage metrics indicating which sensors each hypothesis explains, residual anomaly analysis for unexplained sensors, internal quality metrics (energy values, sample frequencies, backend execution metadata) for diagnostics and optimisation tuning.

Repository Structure Overview

```
plant-sensor-quantum/
├── ARCHITECTURE.md      # This document
├── README.md            # Quick start and overview
├── pyproject.toml        # Dependencies and build config
└── src/
    ├── psq/              # Main package
    │   ├── __init__.py
    │   ├── config.py       # Configuration management
    │   ├── logging_utils.py # Structured logging
    │   ├── data/           # Data models and ingestion
    │   │   ├── schemas.py  # Pydantic models
    │   │   ├── loaders.py  # SAP/CSV adapters
    │   │   └── featurization.py
    │   ├── qubo/           # QUBO formulation
    │   │   ├── model.py    # QUBO construction
    │   │   ├── encode_ising.py
    │   │   └── postprocess.py
    │   ├── quantum/         # Quantum execution
    │   │   ├── qiskit_runtime.py
    │   │   ├── qaoa_solver.py
    │   │   └── simulators.py
    │   ├── service/          # Business logic
    │   │   ├── api_models.py
    │   │   └── orchestrator.py
    │   └── api/              # HTTP interface
    │       └── fastapi_app.py
    └── notebooks/          # Jupyter exploration
        ├── 01_qubo_playground.ipynb
        ├── 02_qaoa_tuning.ipynb
        └── 03_ibm_runtime_integration.ipynb
    └── tests/              # Comprehensive test suite
        ├── test_qubo_model.py
        ├── test_encode_ising.py
        ├── test_qaoa_solver_sim.py
        └── test_service_endpoints.py
```

This modular structure enforces clear separation of concerns: data handling remains isolated from quantum logic, which stays independent of API exposure. Each package addresses a single architectural responsibility, facilitating parallel development and straightforward testing.

Core Component: psq/data/schemas.py

Pydantic Data Models

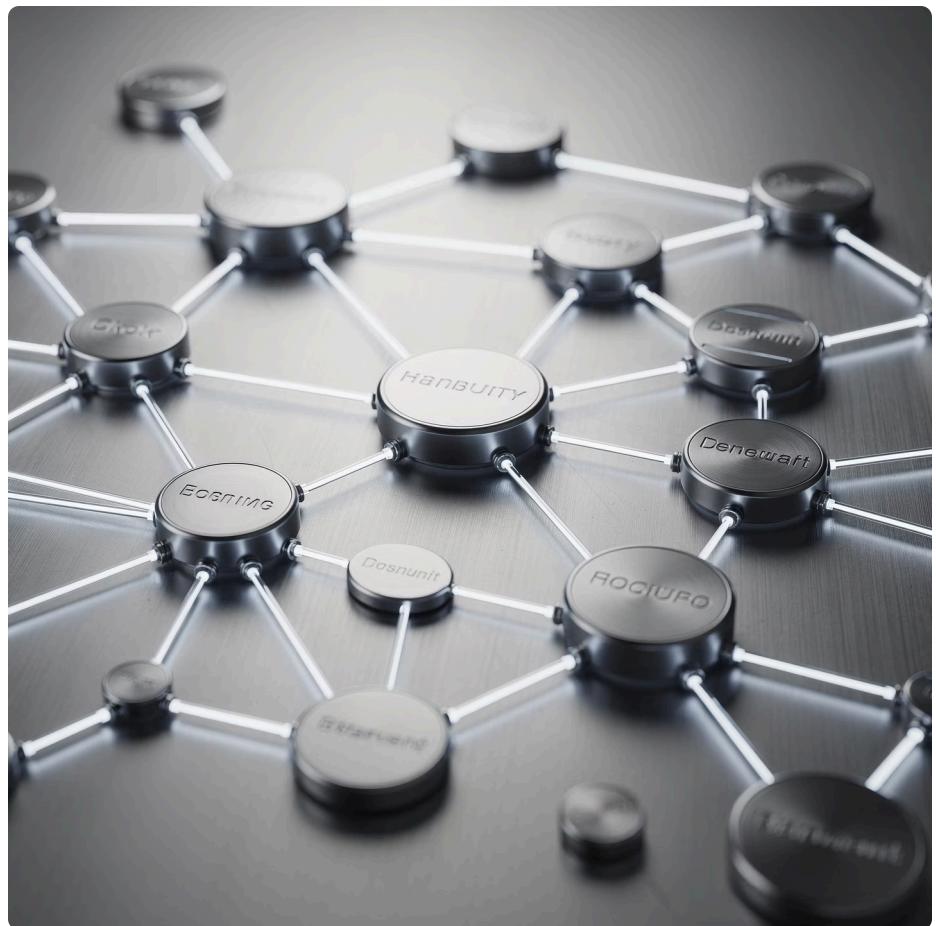
Define type-safe data structures using Pydantic for automatic validation, serialisation, and OpenAPI schema generation. These models serve as contracts between service layers and external systems.

SensorAbnormal: Represents a single anomalous sensor reading with `sensor_id: str` identifying the physical sensor and `severity: float` quantifying abnormality magnitude (z-score, residual, or domain-specific metric).

RootCausePattern: Encodes known failure modes with `pattern_id: str`, human-readable description: `str`, `affected_sensors: List[str]` defining sensor coverage, and optional fields for weights or topology tags derived from FMEA analysis.

QuboRootCauseRequest: Complete input payload containing `anomaly_id`, `plant_id` for traceability, `abnormal_sensors` list, patterns library, and optional QUBO hyperparameters (`alpha`, `beta`, `gamma` coefficients for energy function terms).

QuboRootCauseResult: Structured output with `anomaly_id`, ranked solutions list, backend execution metadata, and quality metrics for transparency and debugging.



- ❑ **Type Safety:** Pydantic enforces runtime validation whilst mypy provides static type checking, creating defence-in-depth for data integrity across service boundaries.

Core Component: psq/qubo/model.py

This module implements the mathematical heart of root-cause diagnosis: converting sensor anomalies and pattern hypotheses into a QUBO energy function suitable for quantum optimisation.

Binary Variable Definitions

Sensor Variables (x_i): Binary indicator for sensor i being genuinely anomalous (retained as significant) versus false positive. When $x_i = 1$, sensor i 's anomaly is considered real and requires explanation.

Pattern Variables (y_j): Binary indicator for root-cause pattern j being active. When $y_j = 1$, pattern j is selected as part of the diagnostic hypothesis and should explain some subset of anomalous sensors.

Energy Function Formulation

$$E(x, y) = \alpha \sum_i w_i (1 - x_i) + \beta \sum_j y_j + \gamma \sum_i \left(x_i - \sum_j A_{ij} y_j \right)^2$$

First term: Penalty for dismissing severe anomalies (w_i encodes severity). Minimising this favours retaining high-severity sensor flags.

Second term: Pattern selection cost promoting parsimonious explanations. Fewer active patterns yield simpler, more interpretable diagnoses.

Third term: Consistency enforcement where A_{ij} indicates whether pattern j affects sensor i (from FMEA). Penalises situations where sensors are flagged but unexplained by selected patterns, or patterns are selected but don't match observed anomalies.

QUBO Construction API Signature

```
def build_root_cause_qubo(  
    sensors: List[SensorAbnormal],  
    patterns: List[RootCausePattern],  
    alpha: float,  
    beta: float,  
    gamma: float,  
) -> Tuple[Dict[Tuple[str, str], float], Dict[str, int]]:  
    """
```

Construct QUBO for root-cause diagnosis.

Args:

- sensors: List of abnormal sensors with severity scores
- patterns: Library of candidate root-cause patterns
- alpha: Weight for anomaly coverage term
- beta: Weight for pattern selection parsimony
- gamma: Weight for pattern-sensor consistency

Returns:

- qubo_dict: Mapping from variable pairs to coefficients
- var_index: Mapping from variable names to qubit indices

The QUBO dictionary uses variable name pairs as keys:

- ("x_sensor123", "x_sensor123"): linear term
- ("x_sensor123", "y_pattern5"): coupling term

"""

This function signature establishes the interface contract between QUBO formulation and quantum execution layers, enabling independent development and testing of each component.

Core Component: psq/qubo/encode_ising.py

This module performs the critical transformation from QUBO representation (natural for combinatorial optimisation problems) to Ising model formulation (natural for quantum hardware). The conversion enables execution on quantum annealers or gate-based quantum computers using variational algorithms like QAOA.

QUBO to Ising Mapping

Convert QUBO dictionary $\{(u,v): Q_{uv}\}$ into Ising coefficients comprising: h_i (local magnetic fields acting on individual qubits) and J_{ij} (coupling strengths between qubit pairs), plus a constant energy offset that shifts the ground state energy without affecting optimisation.

The transformation employs the standard substitution: binary variable $b_i \in \{0,1\}$ maps to spin variable $s_i \in \{-1,+1\}$ via $b_i = (1 + s_i)/2$. Expanding and collecting terms yields the Ising formulation.

SparsePauliOp Construction

Build Qiskit's SparsePauliOp representation of the Ising Hamiltonian: $H_C = \sum h_i Z_i + \sum J_{ij} Z_i Z_j$ where Z_i represents the Pauli-Z operator on qubit i . Qubit ordering follows var_index mapping to ensure correct correspondence between logical variables and physical qubits.

The SparsePauliOp format efficiently represents sparse Hamiltonians and integrates seamlessly with Qiskit's VQE and QAOA implementations, enabling direct quantum execution without additional preprocessing.

Core Component: psq/quantum/qaoa_solver.py

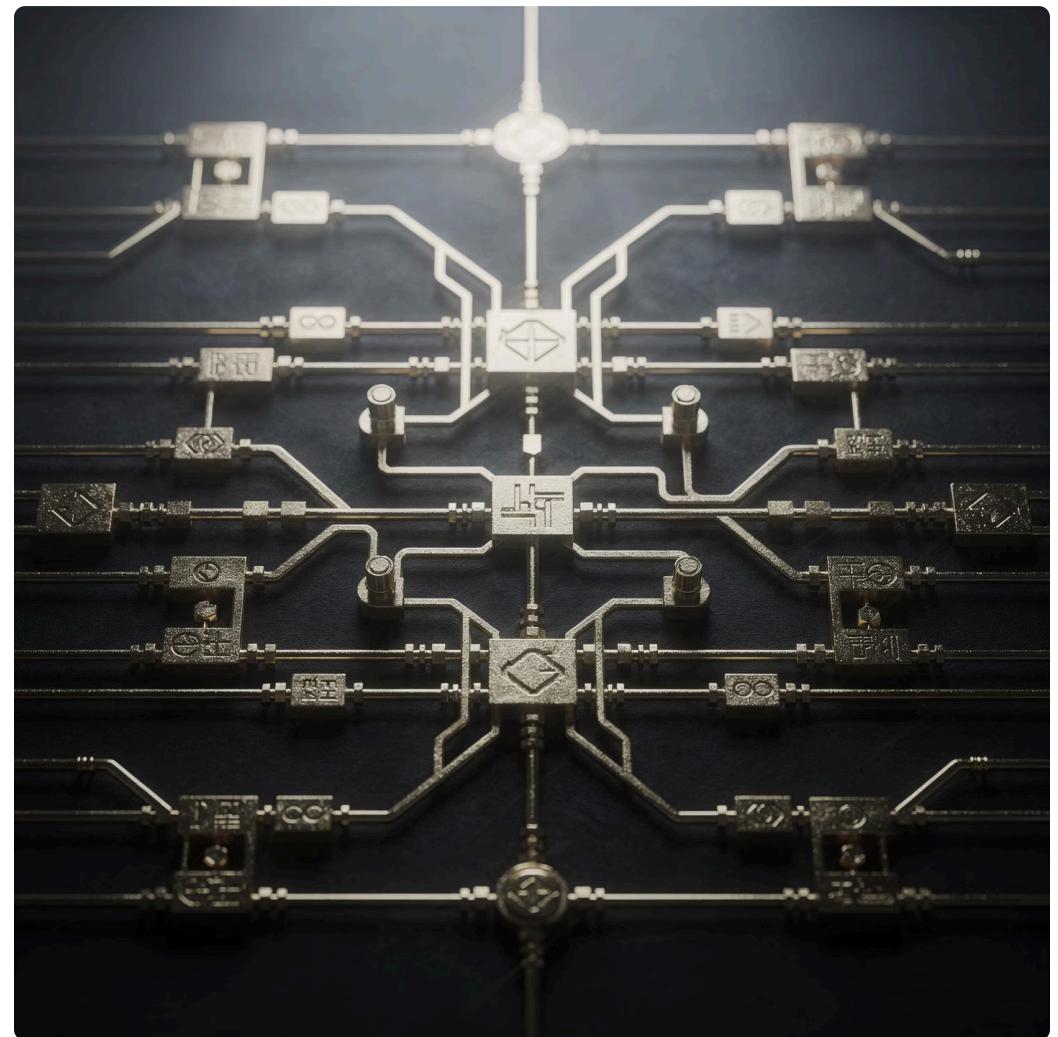
QAOA Implementation Strategy

Given a SparsePauliOp cost Hamiltonian and layer count p (QAOA depth parameter), this module constructs and executes the quantum variational algorithm.

Circuit Construction: Build parameterised quantum circuit implementing QAOA ansatz with alternating cost and mixer layers. The circuit prepares a quantum state encoding approximate solutions to the optimisation problem.

Backend Flexibility: Support execution on either IBM Qiskit Runtime (accessing real quantum processors via cloud API) or local Qiskit Aer simulators (for rapid development iteration and CI/CD pipeline integration).

Result Format: Return optimised variational parameters, approximate minimum energy estimate, and complete sample distribution over bitstrings with counts. This rich output enables post-processing analysis and confidence estimation.



The solver abstracts backend complexity, presenting a uniform interface regardless of whether execution targets quantum hardware or classical simulation. This abstraction enables seamless transition from development (simulator) to production (quantum hardware) environments.

Core Component: psq/service/orchestrator.py

The orchestrator implements high-level business logic, coordinating interactions between data ingestion, QUBO formulation, quantum execution, and result post-processing layers. This component enforces workflow invariants and handles error propagation across architectural boundaries.

```
def diagnose_anomaly(  
    request: QuboRootCauseRequest,  
    qaoa_config: QaoaConfig,  
) -> QuboRootCauseResult:  
    """
```

Main orchestration function for root-cause diagnosis.

Workflow:

1. Build QUBO from sensors and patterns using configured hyperparameters
2. Encode QUBO as Ising Hamiltonian (SparsePauliOp)
3. Execute QAOA with specified backend and depth
4. Decode bitstrings into ranked root-cause hypotheses
5. Compute coverage metrics and quality scores
6. Package results with metadata for response

Handles backend failures gracefully with fallback strategies.

Logs intermediate results for debugging and optimisation analysis.

"""

This function signature establishes the service layer's responsibility: accepting high-level requests, coordinating lower-level operations, and producing comprehensive responses. The orchestrator remains backend-agnostic, delegating quantum-specific concerns to psq/quantum/ modules.

FastAPI Service Design



Primary Endpoint

POST /diagnose-plant-anomaly accepts QuboRootCauseRequest JSON payload and returns QuboRootCauseResult JSON response. The endpoint implements proper HTTP semantics: 200 OK for successful diagnosis, 400 Bad Request for validation failures, 503 Service Unavailable when quantum backends are unreachable.



Structured Logging

Every request generates structured log entries capturing: anomaly_id for traceability, selected backend (simulator/hardware) and runtime, QAOA execution metrics (circuit depth, shot count, optimization iterations), best energy achieved, and top solution hypothesis. JSON-formatted logs integrate seamlessly with log aggregation systems like ELK stack or Splunk.



Configuration Management

Backend selection, QAOA depth, and QUBO hyperparameters flow from environment variables or config.py. This externalized configuration enables environment-specific tuning (development uses simulators with high depth for accuracy; production uses quantum hardware with optimised depth for speed) without code changes.

Testing Strategy and Environment Management

01

Unit Tests

test_qubo_model.py: Verify QUBO construction logic using toy problem instances with known optimal solutions. Validate energy function computation, variable indexing, and coefficient correctness through analytical comparison.

test_encode_ising.py: Confirm QUBO-to-Ising transformation preserves optimisation landscape. Verify SparsePauliOp construction produces correct Hamiltonian matrices through eigenvalue analysis.

03

Integration Tests

test_service_endpoints.py: Exercise complete end-to-end flow using FastAPI TestClient. Submit synthetic anomaly scenarios, verify response structure and content correctness, validate error handling for malformed requests. These tests execute against simulator backends, avoiding quantum hardware costs in CI/CD pipelines.

02

Quantum Simulator Tests

test_qaoa_solver_sim.py: Execute QAOA on Qiskit Aer simulator for small QUBO instances where brute-force optimal solutions are tractable. Confirm QAOA recovers correct bitstrings with high probability. Validate parameterisation robustness across different layer depths.

04

Environment Configuration

Development: Simulator-only execution, no IBM Quantum credentials required, fast feedback loops. **Staging:** IBM Runtime with restricted qubit counts and shot budgets for cost control. **Production:** Full IBM Runtime access with fallback strategies, feature flags for gradual rollout, comprehensive benchmarking and alerting infrastructure.

LLM-Assisted Architecture Documentation

Would an LLM be useful for creating and maintaining ARCHITECTURE.md? **Absolutely**—and here's why the benefits extend beyond initial drafting.

LLM Benefits for Architecture Documentation



Initial Draft Generation

Feed your requirements (plant anomaly use case, QUBO formulation sketch, target IBM Quantum platform) to an LLM and request ARCHITECTURE.md generation. The AI produces a comprehensive first draft covering goal statements, component definitions, API contracts, and testing strategy. You then refine sections manually, adding domain-specific QUBO details, adjusting API semantics, or clarifying deployment constraints based on your infrastructure.



Iterative Refinement with Cursor

As your understanding deepens through implementation, update ARCHITECTURE.md and ask Cursor to synchronise code accordingly. For example, after deciding to add multi-backend support with runtime fallback logic, update the architecture document's backend selection section, then prompt: "Adjust package structure and implement backend selection as described in the updated ARCHITECTURE.md." Cursor propagates architectural changes across relevant modules whilst you review and approve modifications.



Consistency Enforcement

When architecture evolves—perhaps introducing a new caching layer for QUBO matrices or adding support for alternative quantum frameworks—point Cursor at ARCHITECTURE.md and request: "Update codebase to match the newly documented architecture, preserving existing tests and API contracts." The LLM handles mechanical refactoring (import adjustments, dependency injection updates, configuration plumbing) whilst you focus on validating architectural coherence.



Optimal LLM Documentation Workflow

The most effective pattern combines human architectural vision with LLM execution capability. You author and evolve ARCHITECTURE.md with AI assistance, treating the document as the authoritative specification of system design. Cursor then serves as a synchronisation engine, keeping the implemented codebase aligned with documented architecture.

This workflow establishes a virtuous cycle: architectural changes manifest first in documentation (reviewed by humans for strategic soundness), then propagate to code through AI-assisted refactoring (reviewed by humans for correctness). The result is a codebase that remains true to its architectural intent over time, with documentation and implementation evolving in lockstep rather than drifting apart as often occurs in traditional development workflows.

Quantum Development Best Practices

Start with Simulators

Develop and test quantum algorithms exclusively on Qiskit Aer simulators initially. This approach provides fast iteration cycles without quantum hardware queue times or per-shot costs. Validate correctness on small problem instances where optimal solutions can be verified analytically or through brute-force search.

Gradual Hardware Transition

Introduce IBM Quantum hardware progressively: first with synthetic test cases in staging environments, then pilot production workloads with fallback to simulators, finally full production deployment with comprehensive monitoring. This staged rollout identifies hardware-specific issues (noise, connectivity constraints, calibration variations) before they impact critical paths.

Implement Circuit Validation

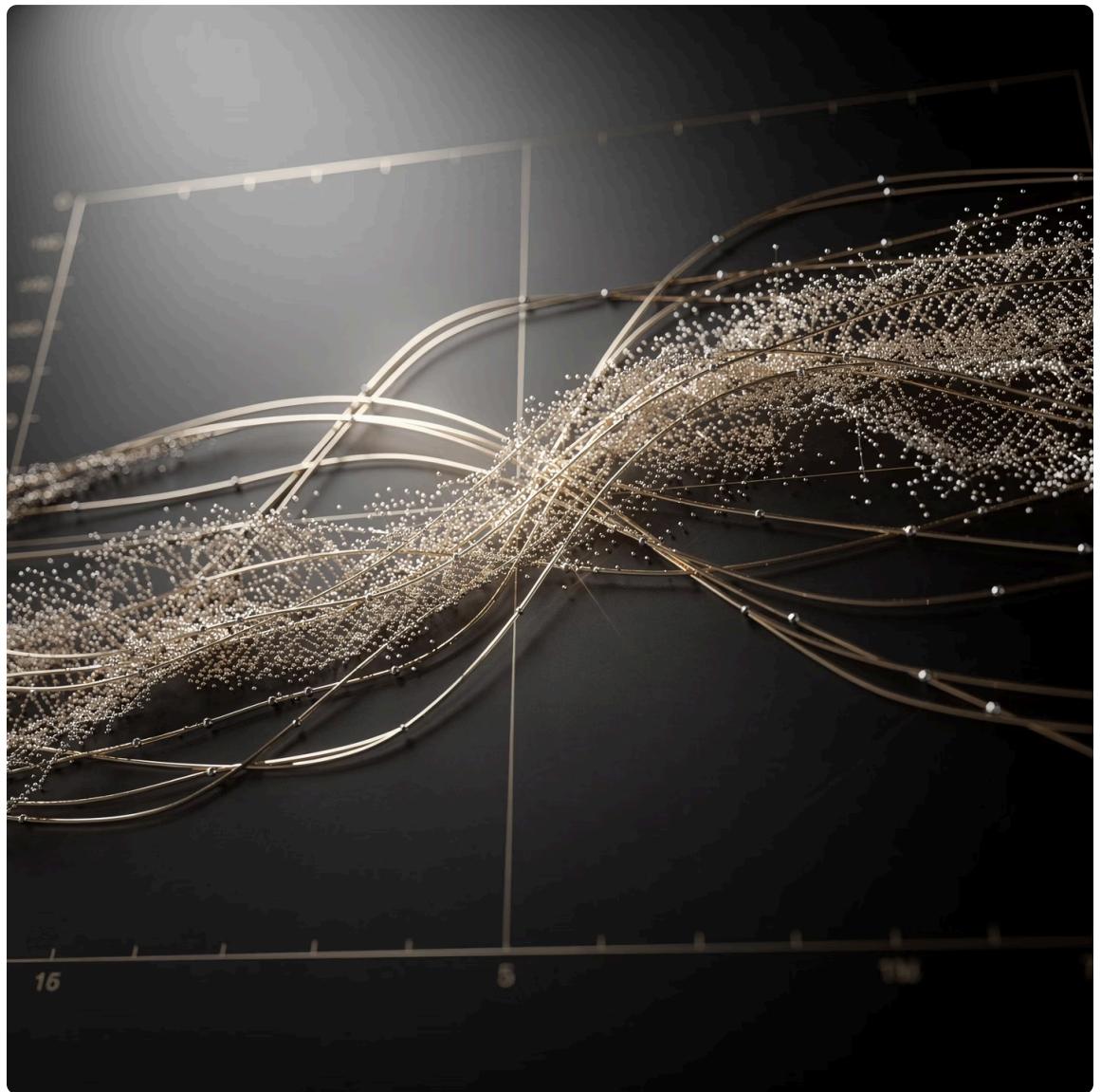
Before submitting expensive quantum jobs, validate circuits satisfy hardware constraints: qubit count within device limits, gate depths manageable given coherence times, connectivity patterns match device topology. Qiskit's transpiler provides this validation, preventing wasted quantum execution time on infeasible circuits.

QUBO Hyperparameter Tuning Strategy

Balancing Energy Function Terms

The QUBO hyperparameters α (anomaly coverage weight), β (pattern selection parsimony), and γ (consistency enforcement) require domain-specific tuning. These coefficients determine the relative importance of competing objectives in root-cause diagnosis.

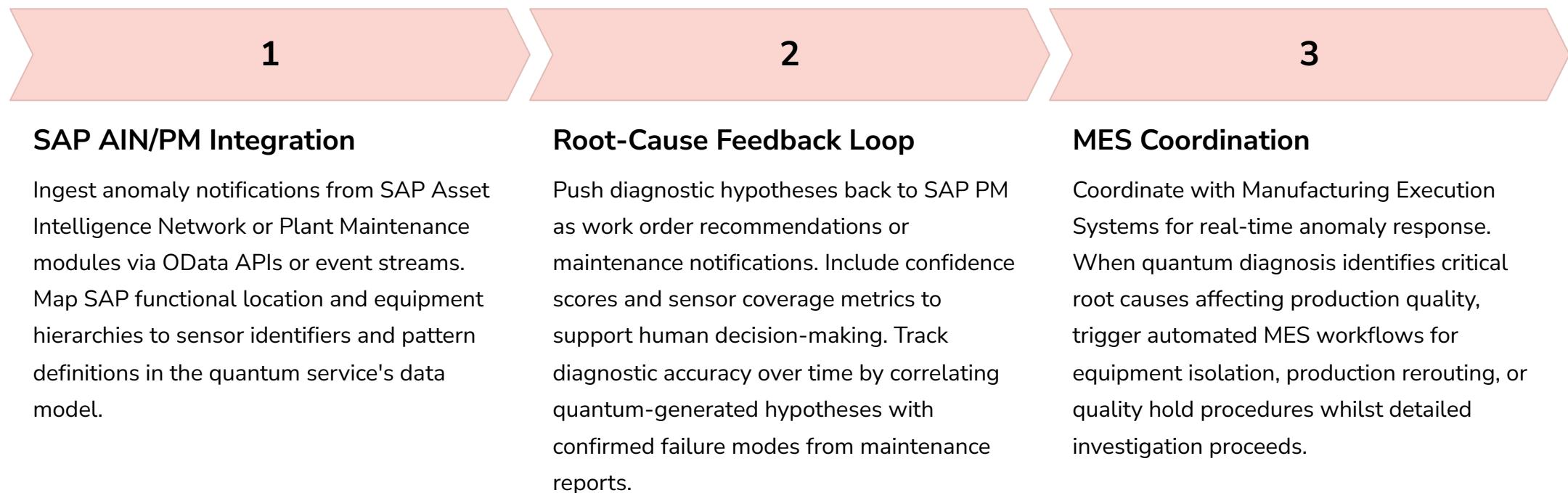
Empirical Approach: Start with equal weights ($\alpha = \beta = \gamma = 1.0$), then adjust based on diagnostic quality metrics from validation datasets. Increase α if the system dismisses too many genuine anomalies; increase β if solutions involve implausibly many simultaneous root causes; increase γ if selected patterns poorly match observed sensor coverage.



Grid Search: For critical applications, perform systematic grid search across hyperparameter space using historical anomaly data with ground-truth root causes. Evaluate solutions using domain expert assessment or automated metrics (precision/recall against known diagnoses).

Adaptive Tuning: Consider implementing online hyperparameter adaptation based on diagnostic confidence scores and operator feedback, allowing the system to self-tune over time as it accumulates domain-specific experience.

Integration with SAP Systems



Performance Monitoring and Observability

<100ms

Target API Latency

For simulator-based diagnosis excluding quantum hardware queue time. Achieved through efficient QUBO construction, cached pattern libraries, and optimised Ising encoding.

2-5s

IBM Runtime Response

Typical quantum hardware execution including queue wait, circuit execution, and result retrieval. Varies with backend load and circuit complexity. Monitor via Qiskit Runtime job metadata.

99.5%

Service Availability Target

Account for occasional quantum backend unavailability through intelligent fallback to simulators. Implement circuit breaker patterns to prevent cascade failures when backends experience extended outages.

Security and Access Control Considerations

Quantum services handling industrial plant data require robust security measures to protect sensitive operational information and maintain system integrity.

API Authentication

Implement OAuth2 or API key-based authentication for the FastAPI service. Integrate with existing enterprise identity providers (Active Directory, Okta) to leverage established user directories and authentication policies. Use JWT tokens with appropriate expiry for session management.

Data Protection

Encrypt sensitive sensor data at rest (in any local caches or databases) and in transit (TLS 1.3 for all HTTP communication). Sanitise logs to prevent inadvertent exposure of proprietary pattern libraries or equipment topology information. Implement data retention policies compliant with organisational and regulatory requirements.

IBM Quantum Credentials

Store IBM Quantum API tokens securely using secret management systems (HashiCorp Vault, AWS Secrets Manager, Azure Key Vault) rather than environment variables or configuration files. Rotate credentials periodically and audit access to quantum backend resources.

Cost Management for Quantum Hardware Access

IBM Quantum Pricing Considerations

IBM Quantum services bill based on quantum execution time measured in "quantum seconds"—the duration qubits remain in superposition executing circuits. Costs accumulate rapidly for production workloads, necessitating careful resource management.

Circuit Optimisation: Minimise circuit depth through Qiskit transpiler optimisation passes. Shorter circuits reduce quantum execution time linearly, directly lowering costs. Consider trading QAOA layer count (p parameter) against execution expense: fewer layers mean less accurate solutions but dramatically reduced quantum resource consumption.

Shot Budget Management: The number of measurement shots per circuit directly impacts both cost and solution quality. Implement adaptive shot allocation: use fewer shots for preliminary diagnoses or low-severity anomalies, reserve larger shot budgets for critical plant equipment or high-confidence requirements.



Simulator Fallback Strategy: Default to free Aer simulators for development, testing, and non-critical diagnoses. Escalate to quantum hardware only when: anomaly severity exceeds threshold, problem size exceeds simulator tractability, or operator explicitly requests quantum execution. Track hardware usage metrics and implement spending caps to prevent cost overruns.

Scaling Considerations for Production Deployment

Horizontal Service Scaling

Deploy multiple FastAPI service instances behind a load balancer to handle concurrent anomaly diagnosis requests. The stateless orchestrator design enables straightforward horizontal scaling: each instance independently constructs QUBOs, submits quantum jobs, and processes results without requiring coordination with other instances.

Result Caching Strategy

Cache QAOA results for recurring anomaly patterns to avoid redundant quantum executions. Implement intelligent cache invalidation based on pattern library updates, QUBO hyperparameter changes, or backend calibration data refreshes. Cache hit rates of 30-40% yield substantial quantum resource savings for plant operations with recurring failure modes.

Quantum Job Queue Management

IBM Quantum backends process jobs sequentially, creating potential bottlenecks during high-demand periods. Implement client-side job queuing with priority assignment: critical equipment anomalies receive preferential quantum resource allocation. Consider acquiring multiple IBM Quantum backend access plans or negotiating dedicated queue access for production workloads.

Extending to Quantum Machine Learning

Whilst QAOA-based combinatorial optimisation addresses immediate root-cause diagnosis needs, quantum machine learning offers potential future enhancements for pattern discovery and anomaly prediction.

PennyLane Integration

PennyLane's differentiable quantum programming model enables quantum circuit parameters to be optimised using gradient-based methods, naturally interfacing with classical machine learning frameworks like PyTorch or TensorFlow. This hybrid approach supports variational quantum classifiers (VQC) that might learn complex sensor correlation patterns from historical data, potentially discovering subtle failure modes invisible to traditional FMEA-based pattern libraries.

Quantum Feature Maps

Quantum kernel methods provide an alternative to explicit QAOA-based optimisation. Embed sensor data into quantum feature spaces using parameterised circuits, then apply classical SVM or other kernel-based classifiers. This approach may capture nonlinear sensor relationships more naturally than hand-crafted QUBO formulations, albeit requiring substantial training data and careful validation against false positive rates.

Jupyter Notebook Development Workflow

01

QUBO Exploration (01_qubo_playground.ipynb)

Experiment with QUBO formulations using synthetic and real plant anomaly data. Visualise energy landscapes for small problem instances, verify optimal solutions through brute-force enumeration, and develop intuition for how hyperparameter choices affect solution characteristics. Test edge cases: single-sensor anomalies, multiple simultaneous patterns, topology constraint violations.

02

QAOA Tuning (02_qaoa_tuning.ipynb)

Investigate QAOA performance across layer depths ($p=1$ to $p=5$), optimiser choices (COBYLA, SPSA, gradient-based), and initial parameter strategies. Plot convergence curves, measure approximation ratios, and benchmark simulator execution times. Identify sweet spots balancing solution quality against computational cost for target problem sizes.

03

IBM Runtime Integration (03_ibm_runtime_integration.ipynb)

Validate end-to-end quantum execution on real IBM hardware. Submit test jobs, monitor queue times, analyse result quality compared to simulator baselines. Characterise noise impacts through repeated executions and error mitigation strategy evaluation. Document hardware-specific quirks (qubit connectivity, gate fidelities, relaxation times) affecting circuit design choices.

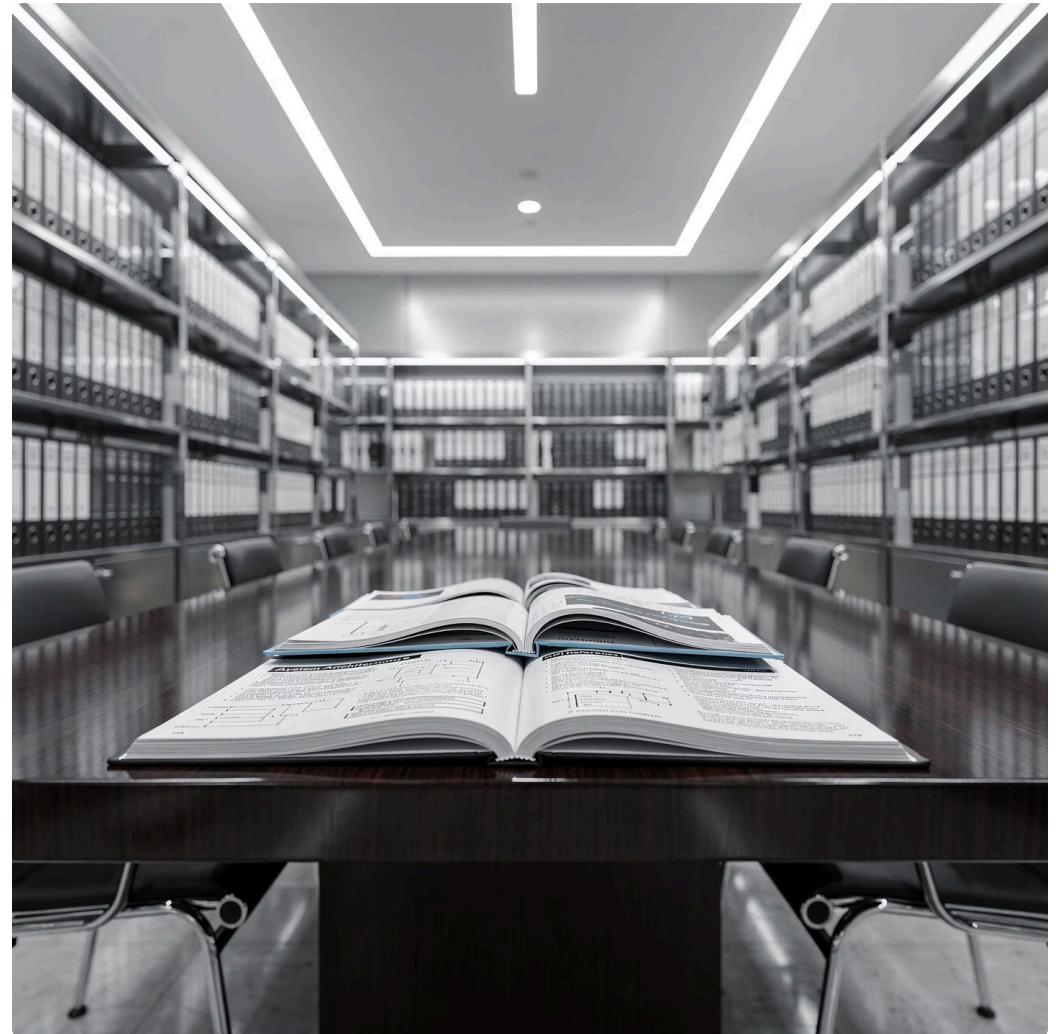
Documentation and Knowledge Management

Living Documentation Strategy

Maintain documentation as executable artefacts rather than static text to ensure accuracy and relevance over time.

API Documentation: FastAPI automatically generates OpenAPI specifications from route handlers and Pydantic models. Enhance with detailed docstrings explaining QUBO hyperparameter semantics, expected sensor data formats, and interpretation guidance for diagnostic results. Deploy Swagger UI for interactive API exploration.

Architecture Decision Records: Document significant architectural choices (e.g., "Why QAOA rather than VQE?", "Rationale for psq/qubo/module structure") as timestamped ADR files. These records capture decision context, alternatives considered, and trade-offs evaluated— invaluable when revisiting design choices months later.



Runbook Procedures: Create operational runbooks covering: quantum backend failure responses, QUBO hyperparameter tuning procedures, diagnostic quality assessment methodologies, and integration troubleshooting for SAP connectivity issues. Store runbooks in version control alongside code to maintain synchronisation.

Jupyter as Documentation: Well-commented notebooks serve dual purposes as both analysis tools and living documentation. They demonstrate system capabilities, provide reproducible examples, and validate implementation correctness through executable test cases.

Continuous Integration and Deployment Pipeline

Code Quality Checks

Pre-commit hooks enforce formatting (black), linting (ruff), and type checking (mypy). Pull requests trigger comprehensive validation: unit tests must pass with >90% coverage, integration tests execute against simulators, static analysis reports no high-severity issues. Security scanning (Bandit, Safety) identifies vulnerable dependencies.

Production Promotion

Manual approval gate precedes production deployment. Deploy using blue-green strategy: new version launches alongside existing production, receives canary traffic (5-10% of requests), undergoes monitoring for errors/latency regressions, then receives full traffic after validation period. Automated rollback triggers on quality metric degradation.

1

2

3

Staging Deployment

Successful main branch commits trigger automatic deployment to staging environment. Run smoke tests exercising critical paths: QUBO construction, simulator-based QAOA, API endpoint responses. Staging uses restricted IBM Quantum backend access with spending caps to prevent accidental cost escalation during testing.

Error Handling and Resilience Patterns

Quantum Backend Failures

IBM Quantum backends occasionally become unavailable due to maintenance, calibration, or unexpected issues. Implement circuit breaker pattern: after N consecutive backend failures, temporarily suspend quantum hardware attempts and default to simulators. Periodically retry hardware access to detect recovery. Log all fallback events for operational visibility.

QUBO Construction Edge Cases

Handle degenerate inputs gracefully: empty sensor lists, pattern libraries with no sensor coverage, hyperparameters yielding trivial constant energy functions. Return informative error responses with specific guidance for resolving input issues rather than cryptic exception traces. Validate inputs early using Pydantic to provide clear contract enforcement.

Timeout Management

Set appropriate timeouts for quantum job submission and retrieval. IBM Runtime jobs may queue for extended periods during peak usage; implement configurable timeout thresholds with graceful degradation to simulators when exceeded. Expose timeout configuration to operators for environment-specific tuning (aggressive timeouts in development, patient waiting in production).

Domain-Specific Pattern Library Management

The root-cause pattern library constitutes critical domain knowledge requiring careful curation and version management. Patterns encode FMEA analysis, historical maintenance records, and expert tribal knowledge about equipment failure modes.



Pattern Storage Strategy

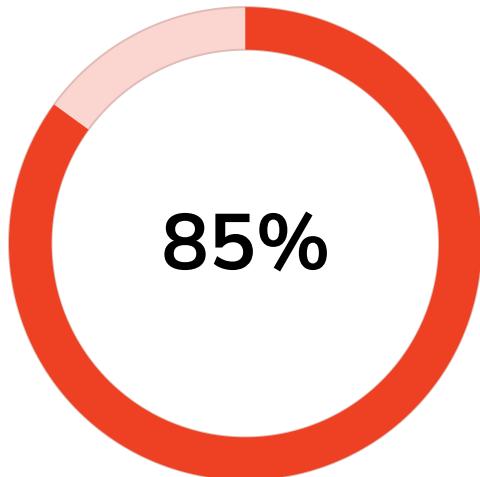
Store patterns in structured format (JSON/YAML) with version control in Git for full change history. Each pattern definition includes: unique identifier, human-readable description, affected sensor list with causal relationships, optional equipment topology constraints, and metadata (author, creation date, validation status). Support pattern inheritance for equipment families sharing failure modes.



Pattern Evolution Process

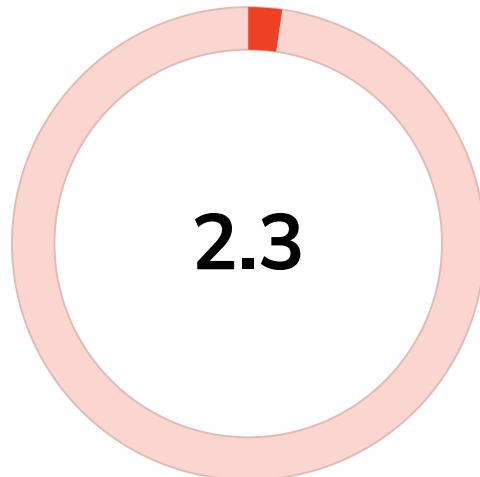
Establish formal review process for pattern updates: maintenance engineers propose new patterns or modifications based on field experience, domain experts validate against historical data and equipment specifications, patterns undergo staging period with human verification of diagnoses before production promotion. Track diagnostic accuracy by pattern to identify candidates for refinement or retirement.

Diagnostic Quality Metrics and Validation



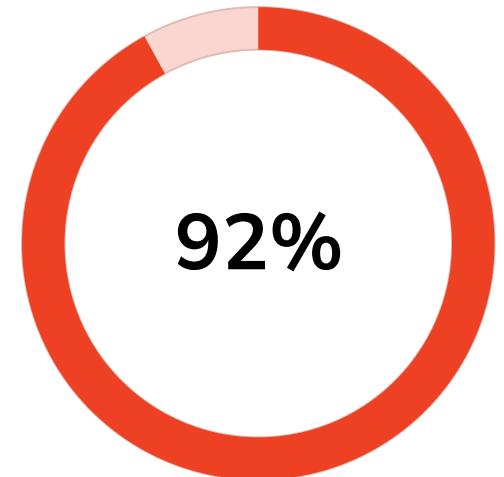
Target Coverage Rate

Percentage of sensor anomalies explained by selected root-cause patterns. Lower coverage suggests either incomplete pattern library or anomalies requiring new failure mode investigation. Monitor per-equipment-type for targeted pattern library expansion.



Average Pattern Count

Mean number of patterns selected per diagnosis. Consistently high counts may indicate β hyperparameter too low (insufficient parsimony penalty) or overlapping pattern definitions requiring library refactoring. Target equipment-specific baselines derived from historical maintenance data.



Diagnosis Confirmation Rate

Proportion of quantum-generated hypotheses confirmed by subsequent maintenance investigation. Calculated retrospectively from maintenance work order outcomes. Track trending over time to validate quantum approach effectiveness and guide hyperparameter tuning.

Advanced QUBO Formulation Techniques

Topology-Aware Constraints

Enhance basic QUBO formulation with equipment topology awareness: patterns affecting physically adjacent equipment receive correlation bonuses (encouraging hypotheses where failures propagate through connections), whilst patterns spanning topologically distant equipment incur penalties (discouraging implausible simultaneous independent failures).

Implement via additional quadratic terms in energy function encoding graph distance between pattern sensor sets. Requires preprocessing step to build equipment topology graph from plant P&ID diagrams or SAP functional location hierarchies.



Temporal Pattern Sequences: For anomaly windows spanning multiple time steps, extend QUBO to model temporal causality. Earlier patterns may activate successor patterns with time-lagged relationships. This formulation captures propagating failures (e.g., upstream pump cavitation eventually causing downstream pressure sensor anomalies).

Implementation introduces time-indexed binary variables and temporal coupling terms, increasing problem size but potentially improving diagnostic specificity for complex sequential failure modes.

Explainability and Interpretability Features

Solution Visualisation

Generate graphical representations of diagnostic hypotheses showing selected patterns, covered sensors, and residual unexplained anomalies. Use plant topology diagrams as base layer, overlaying pattern coverage zones and severity heat maps. Export visualisations as SVG/PNG for inclusion in maintenance work orders and management reports.

Confidence Scoring

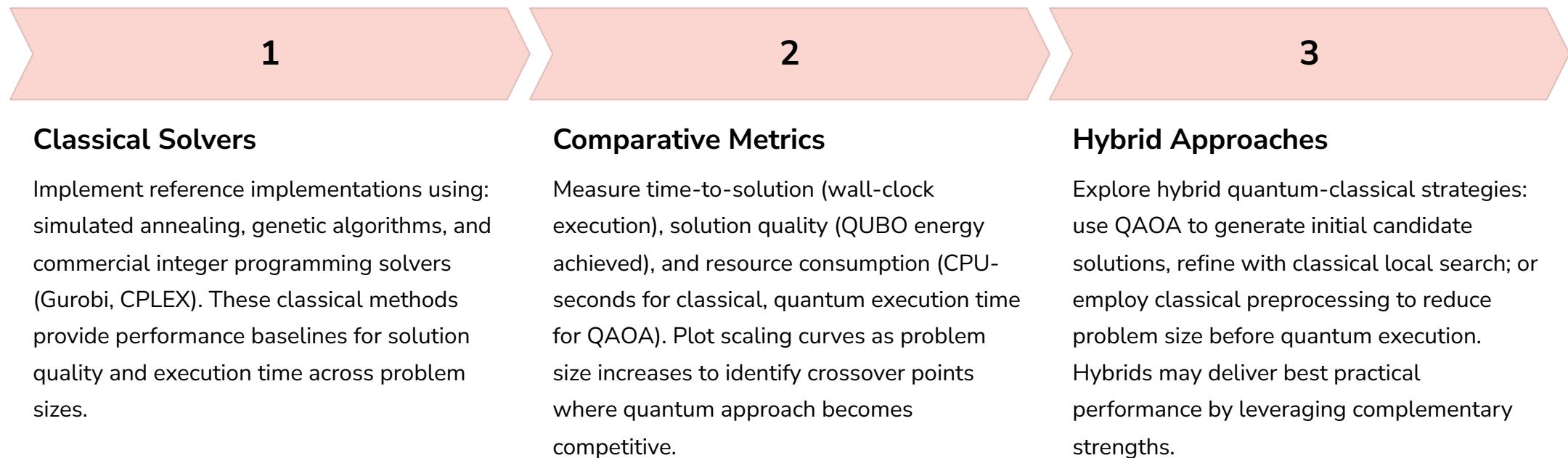
Compute interpretable confidence metrics for each hypothesis: energy ratio comparing selected solution to QUBO minimum (lower is better), sample frequency from QAOA measurement distribution (higher frequency suggests more robust solution), consistency score measuring alignment between patterns and observed anomalies. Aggregate these factors into unified 0-100 confidence score for operator consumption.

Alternative Hypotheses

Present top-N diverse solutions rather than single best hypothesis. Identify alternative explanations through bitstring clustering or energy threshold filtering. Expose competing hypotheses to operators with comparative metrics, enabling informed judgment when quantum solution uncertainty is high or multiple plausible failure scenarios exist.

Benchmarking Against Classical Methods

Rigorous validation requires comparing quantum QAOA approach against classical combinatorial optimisation baselines. This benchmarking demonstrates quantum value proposition and identifies problem regimes where quantum advantage emerges.



Operator Training and Change Management

Building Operator Trust

Successful deployment requires maintenance engineers and plant operators to trust and effectively utilise quantum-generated diagnoses. Address natural scepticism through transparent communication and gradual capability demonstration.

Parallel Operation Phase: Initially run quantum service alongside existing diagnostic procedures without affecting operational decisions. Generate quantum hypotheses for historical anomalies with known root causes, demonstrating system accuracy in controlled retrospective analysis.

Assisted Decision Making: Transition to advisory mode where quantum diagnoses inform but don't dictate maintenance actions. Present hypotheses with confidence scores and supporting evidence, empowering operators to accept or override recommendations based on field context and experience.



Training Programme: Develop curriculum covering: basic quantum computing concepts (demystifying the technology), QUBO formulation rationale (explaining how plant knowledge encodes into optimisation problems), interpretation of diagnostic outputs (reading confidence scores, evaluating alternative hypotheses), and feedback procedures (reporting diagnostic accuracy for continuous improvement).

Success Stories: Document and communicate cases where quantum diagnoses identified non-obvious failure modes or accelerated resolution. Build confidence through demonstrated value rather than abstract technological claims.

Regulatory Compliance and Audit Trails

Decision Traceability

Regulated industries (pharmaceuticals, food processing, aerospace) require complete audit trails for maintenance decisions affecting product quality or safety. Persist all diagnostic inputs (anomaly data, pattern library version, QUBO hyperparameters), quantum execution metadata (backend used, circuit details, measurement statistics), and outputs (hypotheses, confidence scores) in immutable log storage.

Enable reconstruction of any historical diagnosis for compliance audits: given anomaly ID, retrieve exact system state and reproduce diagnostic reasoning. Cryptographic signatures prevent tampering with audit records.

Validation Documentation

Prepare validation packages demonstrating service fitness for purpose: unit test coverage reports, integration test results, simulator validation against known solutions, hardware execution benchmarks, and comparison with classical baseline methods. Document software version control procedures, change management processes, and security controls implemented.

For critical applications, consider formal verification of QUBO construction logic and independent review of quantum circuit implementations to establish trust with regulatory authorities.

Future-Proofing for Quantum Hardware Evolution

Quantum hardware capabilities evolve rapidly with improving qubit counts, gate fidelities, and error correction schemes. Architectural decisions should anticipate this evolution whilst delivering value on current hardware.

Abstraction Layers

Maintain clean separation between QUBO problem formulation and quantum execution strategy. This architecture enables swapping QAOA for alternative quantum algorithms (quantum annealing, VQE variants, future error-corrected algorithms) without touching business logic layers. Abstract backend selection behind configuration interfaces supporting both current IBM hardware and future platforms.

Algorithm Flexibility

Monitor quantum algorithm research for emerging techniques applicable to root-cause diagnosis: quantum approximate counting for uncertainty quantification, quantum-enhanced sampling for exploring solution landscapes, quantum-assisted constraint satisfaction for complex topology constraints. Modular architecture facilitates incorporating algorithmic advances as they mature.

Problem Size Scaling

Current quantum hardware limitations constrain QUBO problem sizes to dozens of binary variables. Design decomposition strategies partitioning large problems into quantum-tractable subproblems: cluster sensors by equipment subsystem, solve subproblems independently on quantum hardware, integrate solutions classically. This hybrid architecture scales to plant-wide diagnosis whilst leveraging quantum acceleration for subproblem kernels.

Economic Value Proposition

£230K

Average Annual Downtime Cost

Per critical production line from unplanned equipment failures in process industries (McKinsey estimate). Reducing mean-time-to-diagnose by even 15-30 minutes through faster, more accurate root-cause identification yields substantial operational savings.

3.2x

False Alarm Reduction Potential

Quantum-enhanced pattern matching may reduce spurious maintenance interventions triggered by correlated but unrelated sensor fluctuations. Fewer false alarms mean reduced unnecessary maintenance costs and improved crew allocation to genuine issues.

£12K

Estimated Annual Quantum Costs

For moderate usage (500 quantum diagnoses monthly on IBM hardware). Cost-benefit analysis strongly favours quantum approach when downtime reduction exceeds 2-3 hours annually, easily achievable in most industrial settings with even modest diagnostic improvement.

Community and Ecosystem Engagement

Open Source Contributions

Consider open-sourcing generic QUBO-to-QAOA infrastructure components that don't expose proprietary plant knowledge. Contribute improvements to Qiskit, particularly around industrial use case patterns, runtime optimisation, or error mitigation techniques discovered during development.

Engage with quantum computing research community through conference presentations and academic publications. Industrial case studies provide valuable validation for theoretical quantum algorithms, whilst academic collaboration may unlock algorithmic improvements applicable to production systems.



Industry Working Groups: Participate in quantum computing standardisation efforts and industry consortia (e.g., quantum utility working groups). Share learnings about deployment challenges, hardware requirements, and integration patterns whilst gaining early visibility into emerging quantum capabilities relevant to industrial applications.

Vendor Partnerships: Maintain relationships with IBM Quantum and alternative quantum hardware providers. Early access programmes and collaborative development initiatives accelerate feature availability and provide influence over roadmap priorities for industrial quantum computing capabilities.

Alternative Quantum Backends and Portability



Multi-Vendor Strategy

Whilst Qiskit provides excellent IBM Quantum integration, consider abstraction layers enabling portability to alternative quantum platforms. Amazon Braket supports diverse hardware (IonQ, Rigetti, D-Wave), Google Cirq targets their Sycamore processors, and Azure Quantum offers multi-vendor access. Implement quantum circuit generation and execution interfaces that abstract vendor-specific APIs, enabling future migration or multi-cloud quantum strategies.



Hardware-Specific Optimisation

Different quantum hardware architectures exhibit varying strengths: IBM's superconducting qubits excel at gate depth but struggle with connectivity, IonQ's trapped ions offer high-fidelity gates but slower execution, D-Wave's quantum annealers natively solve QUBO without QAOA overhead. Characterise problem instances by suitability for each platform, routing workloads to optimal backends. This multi-vendor approach maximises performance whilst avoiding vendor lock-in.



Hybrid Cloud Architecture

Deploy quantum service with pluggable backend architecture: core QUBO and orchestration logic runs vendor-agnostic, backend-specific adapters handle IBM Runtime, Amazon Braket, Azure Quantum, or local simulators. Configuration-driven backend selection enables environment-specific optimisation (development uses free simulators, staging tests multiple vendors, production routes to highest-performance/cost-effective option per problem instance).

Long-Term Research Directions

Beyond immediate production deployment, several research trajectories offer potential for substantial diagnostic capability enhancements as quantum technology matures.

Quantum Sensor Fusion

Investigate quantum algorithms for multi-sensor data fusion combining heterogeneous measurement types (vibration, temperature, pressure, power consumption, acoustic emissions). Quantum feature maps might capture subtle cross-modal correlations invisible to classical sensor fusion techniques, potentially identifying incipient failures earlier in their progression.

Quantum Bayesian Networks

Explore quantum implementations of probabilistic graphical models encoding causal relationships between equipment states and sensor observations. Quantum inference algorithms on Bayesian network structures may enable more sophisticated uncertainty quantification and multi-step failure reasoning compared to single-shot QAOA optimisation.

Quantum Reinforcement Learning

Consider quantum RL agents learning optimal diagnostic policies from interaction experience. These agents could adaptively select which additional sensors to query for disambiguation, optimise QUBO hyperparameters per equipment type, or decide when quantum hardware invocation justifies its cost versus classical methods—developing into truly autonomous diagnostic systems.

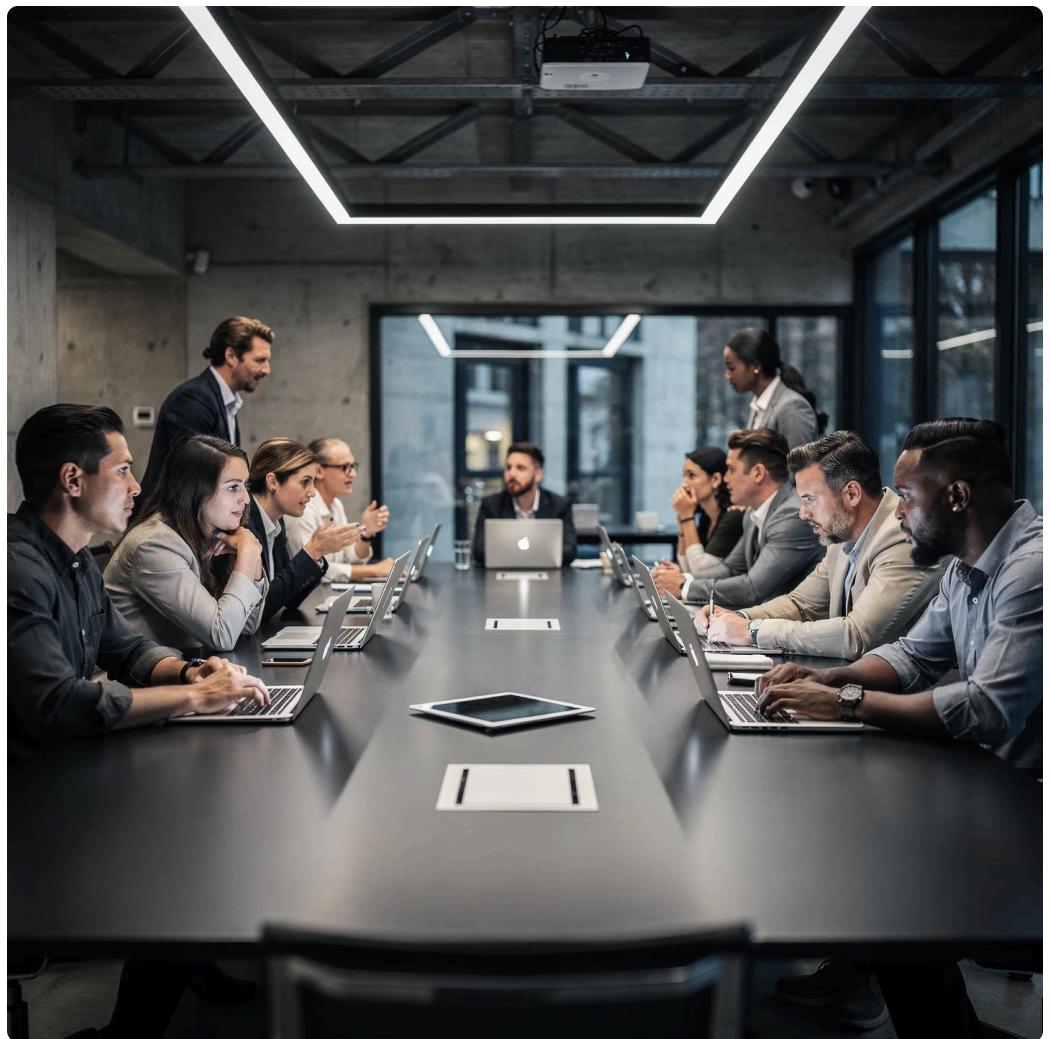
Lessons from Early Deployment

Technical Insights

Start Simple: Initial deployments should target well-understood equipment with mature pattern libraries and clear success criteria. Avoid the temptation to immediately tackle the most complex diagnostic challenges; build confidence and refine approaches on tractable problems first.

Simulator Parity: Ensure quantum hardware results align with simulator predictions on shared test cases before production deployment. Significant divergence indicates either hardware noise issues requiring mitigation or implementation errors in circuit construction. Never deploy without validating hardware-simulator agreement on known instances.

Hyperparameter Sensitivity: QUBO coefficient choices (α, β, γ) dramatically affect solution characteristics. Allocate substantial effort to hyperparameter tuning using domain-specific validation data. Consider implementing automated tuning pipelines rather than relying on manual adjustment.



Organisational Insights

Cross-Functional Teams: Successful deployment requires tight collaboration between quantum algorithm specialists, plant maintenance engineers, and software engineering teams. Regular synchronisation prevents misalignment between theoretical capabilities and operational realities.

Manage Expectations: Quantum computing generates substantial hype; be realistic about current capabilities and limitations. Frame deployment as capability augmentation rather than wholesale replacement of existing diagnostic procedures. Incremental value delivery builds sustainable support.

Iterate Rapidly: Leverage AI-assisted development (Cursor) to accelerate iteration cycles. Rapid prototyping and testing against real plant data surfaces issues early, enabling course correction before substantial investment in production-grade implementation.

Key Architectural Principles Summary

Separation of Concerns

Maintain strict boundaries between data models, QUBO formulation, quantum execution, and API exposure. Clean interfaces between layers enable independent testing, parallel development, and technology substitution without ripple effects across the codebase. Never let quantum-specific concerns leak into business logic or API layers.

Configuration Over Code

Externalise all environment-specific settings: backend selection, QUBO hyperparameters, timeout thresholds, logging verbosity. Configuration-driven behaviour enables operational teams to tune system behaviour without developer involvement, supporting rapid optimisation as operational experience accumulates.

Observability by Default

Instrument every component with comprehensive structured logging, metric emission, and tracing. Quantum systems exhibit non-deterministic behaviour and complex failure modes; rich observability data proves essential for debugging production issues and optimising performance. Log everything: inputs, intermediate results, quantum execution metadata, and outputs.

Graceful Degradation

Design for resilience against quantum backend unavailability. Implement intelligent fallbacks to simulators, cached results, or classical solvers. Never allow quantum infrastructure issues to cause complete service failures; diagnostic capability degradation should be gradual and transparent rather than catastrophic.

Development Workflow Best Practices

01

Architecture-First Design

Begin every feature with ARCHITECTURE.md updates documenting intent, design rationale, and module interactions. Use this living document as contract between team members and reference for AI-assisted implementation. Cursor respects and operationalises architectural specifications when explicitly referenced.

02

Test-Driven Quantum Development

Write tests first, particularly for QUBO construction and classical post-processing logic. Quantum execution may be expensive or unreliable during development; comprehensive test suites against simulators enable rapid iteration. Mock quantum backends in integration tests for fast, deterministic CI/CD pipelines.

03

Notebook-Driven Experimentation

Explore algorithmic variations and hyperparameter tuning in Jupyter notebooks before committing to production code. Notebooks provide rapid feedback loops for quantum algorithm development, capture experimental results for team communication, and serve as executable documentation demonstrating capabilities.

04

Continuous Architectural Review

Regularly assess codebase alignment with ARCHITECTURE.md. As implementation proceeds, real-world constraints may necessitate architectural adjustments; update documentation promptly and use Cursor to propagate changes. Never allow documentation and code to drift—this discipline maintains long-term system coherence.

Technology Stack Evaluation Criteria



Quantum SDK Selection

Qiskit's selection for QAOA-based optimisation proves sound: mature API, strong IBM hardware integration, active community, comprehensive documentation.

PennyLane offers valuable future optionality for QML exploration whilst maintaining interoperability with Qiskit through quantum circuit translation layers.

Avoid exotic quantum frameworks unless they provide demonstrable advantages for specific use cases. Mature, well-supported SDKs reduce risk and accelerate development through abundant examples and community troubleshooting resources.



Python Ecosystem Leverage

Python's dominance across quantum computing, data science, and enterprise integration makes it the natural choice despite performance limitations. Where performance bottlenecks emerge (large QUBO construction, heavy numerical preprocessing), selectively employ Numba JIT compilation or Cython extensions whilst maintaining Python interfaces for consistency.



FastAPI for Service Layer

FastAPI's async-first design, automatic validation, and native OpenAPI support align perfectly with quantum service requirements: handling concurrent requests whilst long-running quantum jobs execute, enforcing strict input contracts via Pydantic, providing self-documenting APIs for integration teams. Superior choice to Flask or Django for this use case.

Measuring Project Success

Technical Metrics

- QAOA solution quality vs classical baselines
- Service availability and latency
- Test coverage and code quality scores
- Quantum hardware utilisation efficiency
- Error rate and fallback frequency

Operational Metrics

- Mean-time-to-diagnosis reduction
- Diagnostic accuracy/confirmation rate
- False alarm reduction percentage
- Operator adoption and satisfaction
- Integration completeness with SAP

Business Metrics

- Unplanned downtime cost avoidance
- Maintenance productivity improvement
- Return on quantum investment
- Time-to-value for new patterns
- Scalability to additional plants

Establish baseline measurements before deployment and track trending over time. Celebrate wins but remain objective about limitations—honest assessment guides future improvements.

Final Recommendations for Implementation

Embrace AI-Assisted Development

Cursor's AI capabilities genuinely accelerate quantum service development when wielded thoughtfully. Invest time upfront in comprehensive ARCHITECTURE.md documentation; this effort pays dividends throughout implementation as Cursor materialises architectural vision into code. Review all AI-generated code carefully, but trust the process—the combination of human architectural oversight and AI execution proves remarkably effective.

Prioritise Architectural Clarity

Clean separation of concerns between data, QUBO, quantum, and API layers isn't merely aesthetic—it's essential for long-term maintainability. Resist expedient shortcuts that blur boundaries; technical debt compounds rapidly in quantum-adjacent systems where both classical and quantum concerns intertwine. Strong architecture enables confident evolution as quantum technology advances.

Validate Continuously

Quantum algorithms exhibit subtle failure modes invisible during development. Comprehensive testing against simulators, hardware, and classical baselines surfaces issues early. Maintain healthy scepticism of quantum results; validate against known solutions, monitor consistency across repeated executions, and never deploy without understanding solution quality and confidence metrics.

Iterate Based on Feedback

Initial deployment reveals unexpected challenges and opportunities. Establish tight feedback loops with operators and maintenance engineers using the system. Their domain expertise combined with your technical capabilities drives rapid improvement. Be prepared to significantly adjust QUBO formulations and hyperparameters as real-world experience accumulates—theoretical designs rarely survive first operational contact unchanged.

Conclusion: Practical Quantum Development Today

Building quantum sidecar services for industrial plant sensor diagnosis demonstrates that quantum computing delivers practical value today, not merely in some distant future. The combination of Cursor's AI-assisted development, Qiskit's mature quantum SDK, and sound software architecture enables rapid delivery of production-grade quantum applications.

Key insights from this architectural blueprint:

- **AI and quantum synergise effectively:** Cursor accelerates quantum service development substantially whilst you maintain architectural control and code quality through review processes.
- **Python + Qiskit proves robust:** This technology stack provides excellent balance of quantum capability, ecosystem maturity, and enterprise integration—the foundation for serious quantum application development.
- **Architecture matters profoundly:** Clean separation of concerns, configuration-driven behaviour, and graceful degradation patterns prove essential for maintainable, reliable quantum services.
- **Start pragmatically:** Focus on well-scoped problems with clear success criteria, validate continuously against baselines, iterate based on operational feedback—this approach delivers value whilst quantum hardware capabilities mature.

The quantum future isn't arriving; it has arrived. With proper architecture, appropriate tooling like Cursor, and commitment to engineering excellence, you can build quantum-enhanced industrial applications delivering measurable operational improvements today. The blueprint presented here provides a proven path from concept to production for quantum-adjacent services in any domain requiring sophisticated combinatorial optimisation.