# Building Resilient Applications: A Deep Dive into Standalone Temporal Python Workers

## 1. Introduction to Temporal Python Workers

In the architecture of a Temporal application, workers serve as the critical bridge between your application code and the Temporal Cluster. They are the execution hosts responsible for polling the Temporal Service for tasks, executing the business logic defined in your workflows and activities, and returning the results back to the cluster. Without workers, the Temporal Cluster, while orchestrating the durable execution, would have no means to actually run your custom application code.

The Temporal Python SDK provides a comprehensive framework for authoring these workflows and activities using the Python programming language. It abstracts away much of the complexity inherent in distributed systems, allowing developers to focus on defining their business logic as durable, fault-tolerant workflows.

Python workers continuously long-poll the Temporal Service for new workflow tasks and activity tasks assigned to their specific task queues. Upon receiving a task, the worker executes the corresponding workflow or activity function and then reports the outcome back to the Temporal Cluster. This polling mechanism and the worker's role in executing code are fundamental to Temporal's ability to ensure that application logic durably executes even in the face of worker failures, network partitions, or other infrastructure issues.

## 2. Setting Up Your Python Development Environment

Before diving into the specifics of defining and running Temporal Python workers, it is essential to establish a suitable development environment. The Temporal Python SDK requires Python version 3.9 or newer to function correctly.

### Environment Setup

A recommended approach for setting up the environment and managing dependencies involves using tools like `uv`. After ensuring Python 3.9+ is installed (e.g., `uv python install 3.13`), developers can clone the Temporal Python samples repository and install all required dependencies by running `uv sync` from the project root.

### Project Structure

For organizing a Temporal application project, a best practice involves structuring the codebase logically. A common and highly recommended project layout includes separate directories for activities, workflows, client code (for starting workflows), and workers.

**Example project structure:**

```
mkdir -p data-processing-service/{activities,workflows,client,workers/python}
```

This separation of concerns enhances code readability, maintainability, and allows for independent scaling and deployment of different components.

## 3. Connecting to the Temporal Server

A standalone Temporal Python worker must establish a connection to a Temporal Cluster to poll for tasks and report results. This connection is managed through a `Client` instance provided by the Temporal Python SDK.

### Client Initialization and Connection Options

The `Client.connect()` method is the primary interface for establishing a connection to the Temporal server. This method typically takes the server's gRPC address as its first argument, such as `"localhost:7233"` for a local development setup. An optional `namespace` argument can also be provided to specify the target Temporal Namespace for operations.

**Important Note**: A `Client` instance in the Python SDK does not have an explicit `close()` method; its lifecycle is typically managed implicitly or through the worker's shutdown.

### Connection Methods

Temporal offers flexibility in how a client connects to its cluster, adapting to various deployment environments:

**Local Development Server**

For rapid local development and initial testing, a Temporal development server can be started using:

```
temporal server start-dev --db-filename temporal.db --ui-port 8080
```

This command automatically sets up a local database, starts the Temporal Web UI (typically at `http://localhost:8080` or `http://localhost:8233`), and creates a default Namespace. Workers connect to this local server using `localhost:7233`.

**Temporal Cloud**

For production use cases or scalable proofs of concept, Temporal Cloud is the recommended environment. Connecting to Temporal Cloud requires specific credentials:

- Temporal Cloud Namespace ID
- Namespace's gRPC endpoint (format: `namespace.unique_id.tmprl.cloud:port`)
- SSL certificate (.pem) and private key (.key) files for mTLS authentication

**Example connection:**

```python
client = await Client.connect(
    "your-namespace.unique-id.tmprl.cloud:7233",
    namespace="your-namespace",
    tls=True,
    tls_cert="path/to/cert.pem",
    tls_key="path/to/key.pem"
)
```

**Self-hosted Temporal Cluster**

For scenarios requiring production-level features with customizability, deploying a self-hosted Temporal Cluster via Docker is an option:

```bash
# Clone and start the cluster
git clone https://github.com/temporalio/docker-compose.git
cd docker-compose
docker compose up
```

Workers need to be configured with the appropriate IP address and port to connect to the Temporal Server container within the Docker network.

## 4. Defining Temporal Workflows in Python

In the Temporal Python SDK, a workflow is defined as a Python class. This class encapsulates the durable, fault-tolerant business logic that orchestrates activities and manages long-running processes.

### Workflow Class Structure

A workflow class is identified by the `@workflow.defn` decorator. This decorator can optionally take a `name` parameter to customize the workflow's registered name, or `dynamic=True` to designate it as a catch-all for otherwise unhandled workflow types.

The core execution logic of a workflow resides in a single method decorated with `@workflow.run`. This method must be an `async def` function, signifying its asynchronous nature. The first parameter of this method must be `self`, followed by any positional arguments that represent the inputs to the workflow.

**Example workflow:**

```python
# workflows.py
from temporalio import workflow
from datetime import timedelta

# Import activity from a separate file
with workflow.unsafe.imports_passed_through():
    from .activities import process_data_activity

@workflow.defn
class DataProcessingWorkflow:
    @workflow.init
    def __init__(self, initial_data: str):
        self.current_data = initial_data
        workflow.logger.info(f"Workflow initialized with: {self.current_data}")

    @workflow.run
    async def run(self, input_data: str) -> str:
        workflow.logger.info(f"Starting DataProcessingWorkflow with input: {input_data}")

        # Execute an activity
        processed_result = await workflow.execute_activity(
            process_data_activity,
            input_data,
            schedule_to_close_timeout=timedelta(seconds=30),
        )

        self.current_data = processed_result
        workflow.logger.info(f"Workflow completed. Final data: {self.current_data}")
        return processed_result

    @workflow.signal
    async def update_data(self, new_data: str):
        self.current_data = new_data
        workflow.logger.info(f"Data updated via signal: {self.current_data}")

    @workflow.query
    def get_current_data(self) -> str:
        return self.current_data

    @workflow.update
    async def append_to_data(self, suffix: str) -> str:
        # Validator for the update
        @append_to_data.validator
        def validate(self, suffix: str):
            if not suffix:
                raise ValueError("Suffix cannot be empty")

        self.current_data += suffix
        workflow.logger.info(f"Data appended via update: {self.current_data}")
        return self.current_data
```

## Deterministic Constraints and the Workflow Sandbox

A cornerstone of Temporal's durability and fault tolerance is the strict requirement that workflow code must be deterministic. This means that given the same history of events, a workflow execution must always produce the same sequence of commands.

**Workflow code must avoid:**

- Non-deterministic functions (random number generators, UUID.randomUUID(), etc.)
- Direct I/O or Service Calls (network I/O, database calls, external services)
- Mutable Global Variables
- Native Threading or Blocking Concurrency
- System Time (use `workflow.now()` instead of `datetime.now()`)

The Temporal Python SDK runs workflow code within a specialized sandbox environment by default. This sandbox restricts access to certain non-deterministic operations, raising errors if violations occur.

## Signal, Query, and Update Handlers

Temporal workflows can interact with external systems and receive external input through signals, queries, and updates:

**Signals (`@workflow.signal`)**

A signal method is used to send asynchronous messages to a running workflow. It can be an `async` or non-`async` method. Signal methods can mutate workflow state and initiate other workflow operations.

**Queries (`@workflow.query`)**

A query method is used to synchronously retrieve the current state of a workflow. It should return a value but must not be an `async` method, nor should it mutate any workflow state.

**Updates (`@workflow.update`)**

Updates are a more recent and robust mechanism for interacting with workflows, offering both input and a return value. They can be `async` or non-`async` and are allowed to mutate workflow state and make calls to other workflow APIs.

# 5. Defining Temporal Activities in Python

Activities in Temporal represent the actual side effects that interact with the outside world. Unlike workflows, activities are not constrained by determinism and can perform any operation, such as database calls, API requests, or complex computations.

## Activity Function Structure

An activity is defined as a Python function decorated with `@activity.defn`. Similar to workflows, a custom name for the activity can be set with a decorator argument, or `dynamic=True` can be used for a catch-all activity.

**Example activities:**

```python
# activities.py
from temporalio import activity
import time
import random
import asyncio

@activity.defn
def process_data_activity(data: str) -> str:
    """
    A synchronous activity that simulates data processing.
    """
    activity.logger.info(f"Processing data: {data}")
    # Simulate a blocking I/O or CPU-intensive task
    time.sleep(random.uniform(1, 3))
    processed_data = data.upper() + "-PROCESSED"
    activity.logger.info(f"Finished processing. Result: {processed_data}")
    return processed_data

@activity.defn
async def fetch_external_resource(url: str) -> str:
    """
    An asynchronous activity that simulates fetching an external resource.
    """
    activity.logger.info(f"Fetching resource from: {url}")
    # Simulate an async network call
    await asyncio.sleep(1)
    return f"Content from {url}"

@activity.defn
def long_running_activity(task_id: str) -> str:
```

```
        """
        A long-running activity that heartbeats and handles cancellation.
        """
        for i in range(10):
            try:
                activity.logger.info(f"Long-running task {task_id}: step {i+1}/10")
                time.sleep(2)  # Simulate work
                activity.heartbeat(f"Progress: {i+1}/10")  # Report progress
            except activity.CancelledError:
                activity.logger.warning(f"Long-running task {task_id} cancelled at step {i+1}")
                raise  # Re-raise to propagate cancellation
        return f"Long-running task {task_id} completed."
```

## Synchronous vs. Asynchronous Activities

Activities can be defined as either synchronous (`def`) or asynchronous (`async def`) functions. Asynchronous activities are generally more performant, especially for I/O-bound tasks, as they allow the worker to handle multiple activities concurrently on a single thread while waiting for I/O operations to complete.

## Heartbeating and Cancellation

For long-running activities, it is crucial to implement heartbeating and handle cancellation:

- **Heartbeating** (`activity.heartbeat()`) informs the Temporal Cluster that the activity is still alive and making progress
- **Cancellation** should be handled gracefully by checking for `activity.CancelledError` and cleaning up resources

## Data Serialization and Converters

All parameters passed to and returned from activity functions must be serializable. The Temporal Python SDK's default data converter supports a wide range of types, including:

- `None`, `bytes`, `google.protobuf.message.Message`
- `dataclasses`, `iterables`, `Pydantic` models
- `IntEnum`/`StrEnum` based enumerates, `UUID`

**Note**: The default converter does not natively support `date`, `time`, or `datetime` objects. For complex or custom data types, developers can implement custom payload converters.

## 6. Registering Workflows and Activities with a Worker

Once workflows and activities are defined, they must be registered with a `Worker` instance. The Worker is the runtime component that polls for tasks, dispatches them to the appropriate workflow or activity function, and manages their execution lifecycle.

**Example worker setup:**

```
# worker_main.py
import asyncio
import concurrent.futures
from temporalio.client import Client
from temporalio.worker import Worker
import logging

# Configure logging for visibility
logging.basicConfig(level=logging.INFO)

# Import your workflow and activities
from .workflows import DataProcessingWorkflow
from .activities import process_data_activity, fetch_external_resource, long_running_activity

async def run_worker():
    # 1. Initialize a Temporal Client
    client = await Client.connect("localhost:7233", namespace="default")
```

```python
    # 2. Configure an activity executor for synchronous activities
    activity_executor = concurrent.futures.ThreadPoolExecutor(max_workers=100)

    # 3. Create a new Worker instance
    worker = Worker(
        client,
        task_queue="my-data-processing-task-queue",
        workflows=[DataProcessingWorkflow],
        activities=[process_data_activity, fetch_external_resource, long_running_activity],
        activity_executor=activity_executor,  # Pass the executor for synchronous activities
        # Other optional worker parameters can be set here
        # max_concurrent_activities=50,
        # max_activities_per_second=100.0,
    )

    logging.info("Worker started, polling task queue 'my-data-processing-task-queue'...")

    # 4. Call run on the Worker to start polling and executing tasks
    try:
        async with worker:
            # Keep the worker running indefinitely until a shutdown signal is received
            await asyncio.Future()  # Awaiting a never-resolving future keeps the worker alive
    except asyncio.CancelledError:
        logging.info("Worker shutdown initiated.")
    finally:
        activity_executor.shutdown(wait=True)
        logging.info("Worker gracefully shut down.")

if __name__ == "__main__":
    asyncio.run(run_worker())
```

## 7. Managing the Worker Lifecycle

The lifecycle of a Temporal Python worker involves starting it to begin polling for tasks and gracefully shutting it down when it is no longer needed.

### Starting and Stopping Workers

A `Worker` instance can be run explicitly using the `worker.run()` method, which is an async function that will continuously poll for tasks until a shutdown is initiated. Alternatively, and often preferred in asynchronous Python applications, the worker can be managed using an async with statement.

**Example of controlled shutdown:**

```python
import asyncio
from temporalio.client import Client
from temporalio.worker import Worker
from .workflows import MyWorkflow
from .activities import my_activity

async def run_worker_with_controlled_shutdown(stop_event: asyncio.Event):
    client = await Client.connect("localhost:7233", namespace="my-namespace")
    worker = Worker(client, task_queue="my-task-queue", workflows=[MyWorkflow], activities=
[my_activity])

    try:
        async with worker:
            await stop_event.wait()  # Worker runs until stop_event is set
            logging.info("Stop event received, initiating worker shutdown.")
    except asyncio.CancelledError:
        logging.info("Worker task was cancelled.")
    finally:
        logging.info("Worker context exited.")
```

```python
    async def main():
        stop_event = asyncio.Event()

        # In a real application, you might set up signal handlers here
        # asyncio.get_event_loop().add_signal_handler(signal.SIGINT, stop_event.set)
        # asyncio.get_event_loop().add_signal_handler(signal.SIGTERM, stop_event.set)

        worker_task = asyncio.create_task(run_worker_with_controlled_shutdown(stop_event))

        # Simulate some work or wait for a condition to stop the worker
        await asyncio.sleep(60)  # Run worker for 60 seconds
        stop_event.set()  # Signal the worker to stop

        await worker_task  # Wait for the worker to complete its shutdown
        logging.info("Application finished.")

    if __name__ == "__main__":
        asyncio.run(main())
```

Graceful Shutdown

For robust deployments, graceful shutdown is paramount. The `Worker` constructor includes a `graceful_shutdown_timeout` parameter, which accepts a `datetime.timedelta`. When this timeout is set, the worker will notify any running activities of an impending graceful shutdown before forcefully cancelling them.

# 8. Concurrency Configuration for Python Workers

Optimizing the concurrency settings of Temporal Python workers is crucial for achieving efficient resource utilization and high throughput.

## Activity Executors

The choice of executor for activities is a key consideration, especially for synchronous (blocking) activities:

### Synchronous Multithreaded Activities

For synchronous activities (functions defined with `def` rather than `async def`), the `activity_executor` worker parameter must be set to an instance of `concurrent.futures.ThreadPoolExecutor`.

### Synchronous Multiprocess/Other Activities

If `activity_executor` is set to an instance of `concurrent.futures.Executor` that is not a `ThreadPoolExecutor` (e.g., `ProcessPoolExecutor`), activities are considered multiprocess or other types of activities.

### Asynchronous Activities

Functions defined with `async def` are asynchronous activities. They are often more performant for I/O-bound operations and do not require any special worker parameters for their execution model within the worker's event loop.

## Worker Concurrency Parameters

The `Worker` constructor provides several parameters to fine-tune concurrency and rate limits:

| Parameter Name | Type | Description | Usage/Implications |
|---|---|---|---|
| `max_concurrent_activities` | `int` | Maximum number of activity tasks that will be given to this worker concurrently. | Directly controls the worker's capacity for parallel activity execution. |

| Parameter Name | Type | Description | Usage/Implications |
|---|---|---|---|
| `max_activities_per_second` | `float` | Limits the number of activities per second that this specific worker will process. | Worker-side rate limiting. Helps manage the load on external resources accessed by activities. |
| `max_task_queue_activities_per_second` | `float` | Sets the maximum number of activities per second that the task queue will dispatch. | Server-side rate limiting. If multiple workers on the same queue have different values set, they will "thrash." |
| `max_cached_workflows` | `int` | Sets the cache size for sticky workflow execution. | Default is 10K. Proper sizing avoids expensive workflow replays when workers die or cache is evicted due to memory pressure. |
| `workflow_task_poller_behavior` | `PollerBehavior` | Specifies the behavior of workflow task polling, including concurrency. | Recommended over deprecated `max_concurrent_workflow_task_polls`. |
| `activity_task_poller_behavior` | `PollerBehavior` | Specifies the behavior of activity task polling, including concurrency. | Recommended over deprecated `max_concurrent_activity_task_polls`. |

## Optimizing Worker Resource Utilization

Effective worker tuning involves balancing resource utilization with responsiveness. A common goal is to utilize approximately 80% of the worker's CPU capacity while ensuring that schedule-to-start latencies remain low.

Properly sizing the `max_cached_workflows` parameter is also vital. This cache stores workflow event histories, allowing workers to avoid replaying the entire history from scratch when processing subsequent workflow tasks for the same workflow.

# 9. Conclusion

Standalone Temporal Python workers are fundamental components for building resilient and scalable distributed applications. This comprehensive guide has detailed the essential aspects of their operation, from connecting to the Temporal Cluster and defining the core business logic in workflows and activities, to managing their lifecycle and configuring concurrency.

The strict determinism requirements for workflow code, enforced by the SDK's sandbox, are not arbitrary limitations but foundational principles that enable Temporal's unique durability and replayability. Similarly, the flexibility to define activities as synchronous or asynchronous, coupled with appropriate executor configurations, empowers developers to integrate diverse computational and I/O-bound tasks effectively.

Effective worker deployment hinges on careful consideration of connection methods, especially the security implications of mTLS for cloud environments. Moreover, meticulous concurrency tuning, including the strategic use of activity executors and precise configuration of worker-side and server-side rate limits, is paramount for optimizing resource utilization and achieving desired throughput.

By embracing these architectural principles and operational best practices, developers can construct robust, fault-tolerant, and highly performant applications on the Temporal platform.