

Universal Stochastic Predictor

Phase 4: IO Layer Initiation

Implementation Team

February 19, 2026

Contents

1 Phase 4: IO Layer Initiation Overview	3
1.1 Scope	3
1.2 Design Principles	3
2 Ingestion and Validation	4
2.1 Implementation Modules	4
2.2 Catastrophic Outlier Filter	4
2.2.1 Implementation Notes	4
2.3 Frozen Signal Alarm	4
2.3.1 Recovery Criteria (V-MAJ-6: Frozen Signal Recovery Ratio)	5
2.4 Staleness Policy (TTL)	6
2.4.1 Implementation Notes	6
3 Telemetry Abstraction	7
3.1 TelemetryBuffer Emission	7
3.1.1 Implementation Notes	7
3.2 No Compute Stalls	7
4 Deterministic Logging	8
4.1 Hash-Based Parity Checks	8
4.2 Audit Hashes	8
5 Snapshot Strategy	9
5.1 Atomic Persistence	9
5.2 Binary Serialization	9
5.2.1 Implementation Notes	9
5.3 Integrity Verification	9
5.4 Atomic Write Protocol	9
5.4.1 Implementation Notes	10
6 Security Policies	11
6.1 Credential Injection	11
6.1.1 Implementation Notes	11
6.2 Version Control Exclusion	11
7 Orchestrator Integration	12
7.1 Ingestion Gate in <code>orchestrate_step()</code>	12
7.1.1 Execution Flow	12
7.1.2 Flag Semantics	12
7.1.3 Early Return on Rejection	13
7.2 PRNG Constant	13
7.3 64-bit Precision Enforcement	13

8 Telemetry Buffer Integration (P2.3)	14
8.1 Motivation	14
8.2 Data Model	14
8.3 Integration into orchestrate_step()	15
8.4 Configuration Parameters	15
8.4.1 Configuration Injection Example	16
8.5 Thread-Safety Model	16
8.6 Backward Compatibility	16
8.7 Usage Example	16
8.8 Benefits	17
9 Compliance Checklist	18
10 Phase 4 Summary	19

Chapter 1

Phase 4: IO Layer Initiation Overview

Phase 4 introduces the asynchronous I/O layer for snapshots, streaming, and telemetry export. The primary design goal is to preserve JAX/XLA throughput by decoupling compute from disk or network latency.

1.1 Scope

Phase 4 covers:

- **Telemetry Buffering:** Non-blocking emission of telemetry snapshots
- **Deterministic Logging:** Hash-based parity checks for CPU/GPU validation
- **Snapshot Strategy:** Atomic persistence of predictor state
- **Ingestion and Validation:** Input filtering, staleness policy, frozen signal detection
- **Security Enforcement:** Credential injection and secret exclusion
- **IO Modules:** validators, loaders, telemetry, snapshots, credentials

1.2 Design Principles

- **No Compute Stalls:** JAX compute threads never block on I/O
- **Determinism:** Logs capture reproducible hashes instead of raw state dumps
- **Security:** No raw signals or secrets in logs
- **Configurability:** Logging intervals and destinations injected via config
- **Integrity:** Snapshots and parity logs are hash-verified

Chapter 2

Ingestion and Validation

2.1 Implementation Modules

Phase 4 IO introduces the following modules:

- `io/validators.py`: Outlier, frozen signal, and staleness checks
- `io/loaders.py`: Ingestion gate and decision flags
- `io/telemetry.py`: Non-blocking telemetry buffer and parity hashes
- `io/snapshots.py`: Binary snapshots, hash verification, atomic writes
- `io/credentials.py`: Environment-based credential injection helpers

2.2 Catastrophic Outlier Filter

Input validation must reject catastrophic outliers when $|y_t| > 20\sigma$ relative to historical normalization. In this case, the system must preserve inertial state and emit a critical alert without advancing the transport update.

- Reject observation and keep current state unchanged.
- Emit a critical alert for audit visibility.
- Do not update JKO/Sinkhorn weights for the rejected step.

2.2.1 Implementation Notes

Outlier detection is implemented as a pure function with configuration-driven thresholds. The ingestion gate returns a decision object that preserves inertial state when an outlier is detected.

2.3 Frozen Signal Alarm

If the exact same value is observed for $N_{freeze} \geq 5$ consecutive steps, emit a `FrozenSignalAlarmEvent`. This invalidates the multifractal spectrum and requires:

- Freeze the topological branch (Kernel D).
- Switch to degraded inference mode.
- Continue monitoring until signal variation resumes.

2.3.1 Recovery Criteria (V-MAJ-6: Frozen Signal Recovery Ratio)

The frozen signal lock is released when variance recovers above a configurable ratio of historical variance for a configurable number of consecutive steps.

Algorithm

$$\text{recovered} = \text{detect_frozen_recovery}(\text{variance_history}, \text{historical_var}, \rho, n_c) \quad (2.1)$$

where:

- `variance_history`: Recent residual variances
- `historical_var`: Baseline variance reference
- $\rho = \text{config.frozen_signal_recovery_ratio}$ (default: 0.1): Recovery threshold multiplier
- $n_c = \text{config.frozen_signal_recovery_steps}$ (default: 2): Confirmation window

Recovery is confirmed when:

$$\text{variance}_t > \rho \cdot \text{historical_var} \quad \text{for } n_c \text{ consecutive steps} \quad (2.2)$$

Implementation

```

1 # In evaluate_ingestion():
2 if frozen:
3     residual_variance = np.var(state.residual_buffer)
4     in_recovery = detect_frozen_recovery(
5         variance_history=[residual_variance],
6         historical_variance=reference_variance,
7         ratio_threshold=config.frozen_signal_recovery_ratio, # V-MAJ-6: Use parameter
8         consecutive_steps=config.frozen_signal_recovery_steps
9     )
10    if in_recovery:
11        frozen = False # Lift the frozen signal flag

```

Configuration Parameters

From PredictorConfig:

Parameter	Default	Purpose
<code>frozen_signal_min_steps</code>	5	Consecutive equal values to trigger alarm
<code>frozen_signal_recovery_ratio</code>	0.1	Variance ratio threshold for recovery (10% of baseline)
<code>frozen_signal_recovery_steps</code>	2	Confirmation window for recovery

Table 2.1: V-MAJ-6 Frozen Signal Recovery Configuration

Benefits

- **Automatic Recovery:** No manual intervention needed when signal variance improves
- **Hysteresis:** Recovery threshold ($\rho = 0.1$) is more lenient than typical alarm threshold, preventing oscillation
- **Configuration-Driven:** All parameters injected from config.toml (zero-heuristics policy)

- **State Preservation:** Maintains frozen flag during low-variance periods, automatically lifts when variance returns
- **Signal Quality Supervision:** Enables secondary observability on signal quality degradation patterns

2.4 Staleness Policy (TTL)

Every observation must carry a timestamp for TTL evaluation. If the target delay exceeds Δ_{max} , the JKO update must be suspended immediately.

- Compute staleness as $\Delta_t = t_{now} - t_{obs}$.
- If $\Delta_t > \Delta_{max}$, skip the transport update.
- Preserve state and record a staleness warning event.

2.4.1 Implementation Notes

Staleness is computed as the difference between current time and observation timestamp. The ingestion decision flags a suspended JKO update when the TTL is exceeded.

Chapter 3

Telemetry Abstraction

3.1 TelemetryBuffer Emission

The JKO orchestrator should emit a `TelemetryBuffer` at the end of each step. This buffer is consumed by a dedicated process outside the JAX execution thread.

- The buffer contains summary metrics (CUSUM, entropy, regime flags, OT cost).
- The compute path only enqueues the buffer and continues.
- The consumer is responsible for serialization and persistence.

3.1.1 Implementation Notes

The telemetry buffer is a bounded, thread-safe queue. Buffer capacity is explicitly injected from `PredictorConfig.telemetry_buffer_capacity` to eliminate implicit defaults (zero-heuristics policy). Parity hashes are emitted on a configurable interval and derived from canonical float64 serialization.

```
1 # Instantiation pattern (capacity injected from config)
2 buffer = TelemetryBuffer(capacity=config.telemetry_buffer_capacity)
```

3.2 No Compute Stalls

JAX compute threads must never block on I/O. Telemetry buffers must be non-blocking and consumed by a separate process or thread outside the JAX execution path.

Chapter 4

Deterministic Logging

4.1 Hash-Based Parity Checks

For hardware parity audits, the logger records SHA-256 hashes of the weight vector ρ and the OT cost at configurable intervals. This permits CPU/GPU parity validation without dumping VRAM data.

- Hash interval configured per deployment.
- Hashes derived from canonical float64 serialization.
- Logs are append-only and immutable.

4.2 Audit Hashes

Parity audits must log SHA-256 hashes of ρ and OT cost at configured intervals. Hash input must be derived from canonical float64 serialization to ensure reproducibility across CPU and GPU.

Chapter 5

Snapshot Strategy

5.1 Atomic Persistence

Snapshots must be persisted atomically to prevent partial writes. The IO layer is responsible for:

- Writing to temporary files and renaming atomically.
- Optional compression configured by policy.
- Coordinating snapshot cadence with telemetry output.

5.2 Binary Serialization

Text formats (JSON, XML) are prohibited for critical snapshots due to latency and ambiguity. Use dense binary formats such as Protocol Buffers or MessagePack.

- Encode all fields deterministically.
- Preserve float64 for numerical fidelity.

5.2.1 Implementation Notes

The snapshot serializer uses MessagePack as the default binary format. Hash verification is performed before state injection.

5.3 Integrity Verification

Each snapshot Σ_t must include a hash footer (SHA-256 or CRC32c). The load routine must verify the hash before injecting state into memory.

- Fail closed if hash verification fails.
- Log integrity failures at critical severity.

5.4 Atomic Write Protocol

To avoid partial writes, persist snapshots to a temporary file and then atomically rename to the target path. The rename step must be the only visible operation to consumers.

- Use a unique temporary filename per snapshot.
- Ensure the target file is replaced atomically.

5.4.1 Implementation Notes

Snapshots are written to a unique temporary file and moved into place using atomic rename. Optional fsync ensures persistence across power loss.

Chapter 6

Security Policies

6.1 Credential Injection

Tokens and API keys must not appear in source code. Credentials must be injected at runtime via environment variables or .env files.

6.1.1 Implementation Notes

Credential helpers read from environment variables or .env files and raise explicit errors on missing values.

6.2 Version Control Exclusion

The repository must exclude .env files and credential directories via .gitignore. Secrets must never be committed.

Chapter 7

Orchestrator Integration

7.1 Ingestion Gate in `orchestrate_step()`

The core orchestration pipeline now integrates ingestion validation as a pre-kernel gate. The `orchestrate_step()` function signature is extended to accept observation metadata:

```
1 def orchestrate_step(
2     signal: Float[Array, "n"],
3     timestamp_ns: int,
4     state: InternalState,
5     config: PredictorConfig,
6     observation: ProcessState,
7     now_ns: int,
8 ) -> OrchestrationResult:
9     """Run a single orchestration step with IO ingestion validation."""

```

7.1.1 Execution Flow

The ingestion gate operates as follows:

1. **Input Validation:** Standard signal length and dtype checks.
2. **Ingestion Decision:** Call `evaluate_ingestion()` with current state, observation, and configuration.
3. **Rejection Logic:** If `accept_observation == False`, reject the entire observation without state update (emergency mode).
4. **Degradation Flags:** Apply `suspend_jko_update` and `freeze_kernel_d` flags to control fusion behavior.
5. **Kernel Execution:** Run kernels A-D; if `freeze_kernel_d == True`, mark kernel D output as frozen.
6. **Fusion Selection:** Skip JKO/Sinkhorn if degraded mode or `suspend_jko_update` is set.
7. **State Update:** Only update InternalState if observation is accepted.

7.1.2 Flag Semantics

The `IngestionDecision` object carries the following flags:

- **accept_observation:** If False, reject and preserve inertial state.

- **suspend_jko_update**: If True, freeze weights and skip Sinkhorn.
- **degraded_mode**: If True, emit degraded inference mode prediction.
- **freeze_kernel_d**: If True, mark kernel D output as frozen (no weight update).
- **staleness_ns**: Staleness in nanoseconds for audit logging.
- **events**: Emitted validation events (outliers, frozen signals, staleness alarms).

7.1.3 Early Return on Rejection

If an observation is rejected (catastrophic outlier), the orchestrator returns a degraded result without advancing the state:

```

1 # If observation is rejected, skip state update entirely
2 if reject_observation:
3     updated_state = state
4 else:
5     updated_state = atomic_state_update(...)
```

7.2 PRNG Constant

To eliminate magic numbers in PRNG splitting, we introduce a module-level constant in `api/prng.py`:

```

1 # api/prng.py: GLOBAL PRNG CONFIGURATION
2 RNG_SPLIT_COUNT = 2 # For kernel execution subkeys
```

This constant is now imported by `core/orchestrator.py` to maintain layer isolation and clarity. All PRNG-related constants reside in the API layer.

7.3 64-bit Precision Enforcement

To ensure Malliavin calculus and Signature computation stability, 64-bit precision must be activated at module import time, before any XLA tracing:

```

1 # api/config.py: JAX CONFIGURATION (at module level)
2 import jax
3 jax.config.update("jax_enable_x64", True)
```

This enforces bit-exact reproducibility across CPU/GPU/FPGA backends and must execute before `ConfigManager` initialization.

Chapter 8

Telemetry Buffer Integration (P2.3)

8.1 Motivation

Phase 2 implementations (P2.1 WTMM, P2.2 SDE stiffness, V-MAJ violations) generate rich diagnostic data during orchestration. To enable post-mortem analysis, compliance audits, and debugging without stalling inference, P2.3 integrates a **non-blocking telemetry buffer** into the orchestration pipeline.

Key requirements:

- **Non-Blocking:** Logging never blocks compute threads (async enqueue only)
- **Thread-Safe:** Multiple consumers can safely drain buffer
- **Audit Trail:** Records capture complete prediction state snapshot
- **Integrity:** Parity hashes verify weights and free energy
- **Config-Driven:** Emission interval and buffer capacity injected from config

8.2 Data Model

Each telemetry record captures:

```
1 @dataclass(frozen=True)
2 class TelemetryRecord:
3     step: int                      # Monotonic counter
4     payload: dict                  # Rich diagnostic data
5
6
7 # Payload structure (P2.3):
8 payload = {
9     "step": 42,
10    "timestamp_ns": 17083080000000000000,   # Nanosecond precision
11    "prediction": 1.234,                      # Fused kernel prediction
12    "weights": [0.25, 0.25, 0.25, 0.25],   # Kernel ensemble weights (rho)
13    "kurtosis": 3.5,                         # Signal regularity
14    "holder_exponent": 0.65,                 # Signal regularity
15    "dgm_entropy": 0.45,                     # Signal regularity
16    "modeCollapse_warning": false,           # V-MAJ-5 flag
17    "degraded_mode": false,                  # V-MAJ-7 hysteresis
18    "emergency_mode": false,                # Circuit breaker
19    "parity_hashes": {
20        "rho_sha256": "abc123...",          # Hash of weight vector
21        "ot_cost_sha256": "def456..."       # Hash of Sinkhorn free energy
22    }
23}
```

8.3 Integration into orchestrate_step()

The orchestrator now accepts optional `telemetry_buffer` and `step_counter` parameters:

```

1 def orchestrate_step(
2     signal: Float[Array, "n"],
3     timestamp_ns: int,
4     state: InternalState,
5     config: PredictorConfig,
6     observation: ProcessState,
7     now_ns: int,
8     telemetry_buffer: Optional[TelemetryBuffer] = None, # P2.3
9     step_counter: int = 0, # P2.3
10 ) -> OrchestrationResult:
11     """
12         Run a single orchestration step with telemetry buffering.
13
14         If telemetry_buffer is None, skips telemetry (backward compatible).
15         If provided, enqueues record when hash_interval triggers.
16     """
17     # ... existing orchestration logic ...
18
19     # P2.3: Telemetry Buffer Integration (before return)
20     if telemetry_buffer is not None:
21         parity_record = parity_hashes(
22             rho=final_rho,
23             ot_cost=float(free_energy) if fusion is not None else 0.0
24         )
25
26         telemetry_payload = {
27             "step": step_counter,
28             "timestamp_ns": timestamp_ns,
29             "prediction": float(fused_prediction),
30             "weights": [float(w) for w in final_rho],
31             "kurtosis": float(updated_state.kurtosis),
32             "holder_exponent": float(updated_state.holder_exponent),
33             "dgm_entropy": float(updated_state.dgm_entropy),
34             "modeCollapseWarning": modeCollapseWarning,
35             "degradedMode": degradedMode,
36             "emergencyMode": emergencyMode,
37             "parity_hashes": parity_record,
38         }
39
40         # Emit only when hash_interval triggers (config-driven)
41         if should_emit_hash(step_counter, config.telemetry_hash_interval_steps):
42             telemetry_record = TelemetryRecord(step=step_counter, payload=
43                 telemetry_payload)
44             telemetry_buffer.enqueue(telemetry_record)
45
46     return OrchestrationResult(...)
```

8.4 Configuration Parameters

Parameter	Default	Purpose
<code>telemetry_buffer_capacity</code>	1024	Maximum records in ring buffer
<code>telemetry_hash_interval_steps</code>	1	Emit telemetry every N steps

Table 8.1: P2.3 Telemetry Configuration

8.4.1 Configuration Injection Example

```
1 # config.toml
2 [io]
3 # Telemetry
4 telemetry_hash_interval_steps = 1          # Emit every step (or 10 for sparser logging)
5 telemetry_buffer_capacity = 1024           # Ring buffer size (zero-heuristics injection)
```

8.5 Thread-Safety Model

TelemetryBuffer uses `threading.Lock` for atomic operations:

- `enqueue()`: Acquires lock, appends record, releases ($O(1)$ amortized)
- `drain()`: Acquires lock, extracts all records, clears buffer, releases
- `size()`: Acquires lock, returns current count, releases

This prevents race conditions when orchestrator (compute thread) enqueues while consumer thread drains.

8.6 Backward Compatibility

P2.3 is fully backward compatible:

- If `telemetry_buffer=None`, no telemetry is emitted (default)
- Existing calls to `orchestrate_step()` without telemetry params continue working
- No changes to compute path (telemetry entirely outside `@jax.jit` scope)

8.7 Usage Example

```
1 from stochastic_predictor.io.telemetry import TelemetryBuffer
2 from stochastic_predictor.api.config import PredictorConfigInjector
3
4 # Initialize
5 config = PredictorConfigInjector().create_config()
6 telemetry_buffer = TelemetryBuffer(capacity=config.telemetry_buffer_capacity)
7
8 # In prediction loop:
9 for step in range(num_steps):
10     result = orchestrate_step(
11         signal=current_signal,
12         timestamp_ns=now_ns(),
13         state=state,
14         config=config,
15         observation=obs,
16         now_ns=now_ns(),
17         telemetry_buffer=telemetry_buffer,      # P2.3: Pass buffer
18         step_counter=step,                      # P2.3: Pass step
19     )
20
21     # Optional: Drain telemetry in background thread
22     if step \% 100 == 0:
23         records = telemetry_buffer.drain()
24         # Write to file/database (non-blocking, doesn't stall orchestrator)
25         write_telemetry_async(records)
```

8.8 Benefits

- **Non-Blocking:** Telemetry enqueue is $O(1)$, never stalls inference
- **Audit Trail:** Complete state snapshots for compliance and debugging
- **Integrity:** Parity hashes enable CPU/GPU parity validation
- **Configurable:** Emission interval and buffer size injected from config
- **Thread-Safe:** Lock-based synchronization for multi-threaded consumers
- **Backward Compatible:** Fully optional, no impact on existing code paths

Chapter 9

Compliance Checklist

- **No Compute Stalls:** All logging is asynchronous
- **Binary Format:** Protocol Buffers or MessagePack for snapshots
- **Atomic Snapshots:** Write-then-rename protocol
- **Deterministic Hashing:** SHA-256 on ρ and OT cost
- **Security:** No raw signals, VRAM dumps, or secrets
- **Integrity:** Snapshot hashes verified before load
- **Config-Driven:** Intervals and destinations are injected
- **Module Coverage:** IO helpers implemented for validation, telemetry, snapshots, and credentials
- **Orchestrator Integration:** IO ingestion gate integrated into `orchestrate_step()`
- **PRNG Constants:** Named constants (`RNG_SPLIT_COUNT`) reside in `api/prng.py`
- **Buffer Capacity Injection:** TelemetryBuffer capacity injected from config (zero-heuristics policy)
- **64-bit Precision:** Enforced at module load time (`api/config.py`) before XLA tracing
- **Layer Isolation:** PRNG constants in API layer, not Core layer

Chapter 10

Phase 4 Summary

Phase 4 introduces a non-blocking I/O architecture that preserves deterministic compute while enabling telemetry, logging, and atomic snapshot persistence.