# Verifying Safe Autonomy: A Cross-Language Framework for Neurorights-Compliant Decision-Making in Bio-Feedback Systems

## Establishing Cross-Language Semantic Parity as a Foundational Invariant

The establishment of cross-language semantic parity is not merely a technical objective but a foundational prerequisite for the integrity, predictability, and safety of the autonomous bio-feedback system. The directive to treat this as a "hard invariant" before proceeding to any other aspect of development underscores a critical understanding of complex distributed systems: if the fundamental decision logic diverges between implementation stacks, all subsequent optimizations and adaptive behaviors become inherently unpredictable and potentially hazardous . The proposed architecture, spanning Rust, JavaScript, and ALN, necessitates a rigorous framework to ensure that regardless of the execution environment, the core logic governing autonomy decisions remains identical. This involves a multi-faceted approach encompassing strict schema enforcement, functional equivalence testing, and the potential application of formal verification techniques. The ultimate goal is to create a system where the semantics of an `AutonomyTraceAttributes` object are interpreted identically across all platforms, from parsing and clamping to the final determination of the `autonomy_level`.

A primary mechanism for achieving this parity is the implementation of a strict, versioned JSON schema, designated as `autonomy-trace.v1.json` . This schema serves as the single source of truth for the structure and constraints of the autonomy trace data format. It must meticulously define every field within `AutonomyTraceAttributes`, including their data types, permissible value ranges, and any inter-field dependencies. For instance, the schema would specify that `stress` and `fatigue` are floating-point numbers within a normalized range, `lifeforce_scalar` is a float between 0.0 and 1.0, and `autonomy_level` is an integer corresponding to predefined policy tiers. By enforcing this schema at every point of serialization and deserialization in Rust, JavaScript, and ALN clients, the system prevents syntactically invalid or semantically ambiguous data

from propagating through the ecosystem . However, simple schema validation is insufficient on its own; it only guarantees structural correctness, not behavioral equivalence. An implementation could parse a valid JSON object differently, perhaps due to variations in how floating-point arithmetic or rounding is handled across different language runtimes and libraries.

To address this, the research plan proposes the creation of a minimal "semantic test harness" . This harness would serve as a robust validation engine to move beyond syntax and test for functional equivalence. The methodology would involve feeding a known set of input data, represented as `AutonomyTraceAttributes`, into both the Rust and JavaScript implementations of the autonomy logic. The expected output—a bit-identical `autonomy_level` and `shot_level_label` decision—would then be asserted. This requires careful preparation of "golden test vectors": a curated collection of representative input states that cover various edge cases, boundary conditions, and typical operating scenarios. These vectors must be generated once and treated as immutable artifacts, ensuring that both implementations are tested against the exact same stimuli. Any divergence in the output would immediately flag a semantic inconsistency, which could stem from differences in the underlying decision algorithms, parsing logic, or even the handling of numerical precision. This continuous integration-style testing provides a practical and repeatable method for maintaining parity as the codebases evolve independently. The inclusion of JVM clients in this schema validation further reinforces the commitment to a consistent data contract across the entire technology stack .

The enforcement of the prescribed decision order—`identity → consent → lifeforce → eco → SCALE/turns`—is central to this invariant . This sequence represents a hierarchy of priorities, and its enforcement must be codified directly into the logic of the autonomy governors in each language. For example, a function responsible for computing the final autonomy level must first check the `identity` and `consent` guards. Only if these initial checks pass would it proceed to evaluate the `lifeforce` and `eco` constraints. Finally, after all constraints are satisfied, the `SCALE/turns` logic would determine the precise level of autonomy to grant. This explicit, step-by-step evaluation order must be mirrored precisely in the Rust, JavaScript, and ALN versions of the governor. Any deviation, such as evaluating `eco` before `lifeforce`, would constitute a violation of the semantic invariant, even if the final decision happened to be the same in some cases. The semantic test harness would be instrumental in catching such logical discrepancies. Furthermore, this rigid ordering ensures that the most critical, non-negotiable constraints are evaluated first, preventing situations where a less severe constraint might inadvertently override a more important one.

Looking toward future-proofing the system, the proposal to mirror temporal-logic properties as runtime assertions around the `CompanionAutonomyGovernor::decide` function points toward a deeper level of assurance . Temporal logic allows for the formal specification of system properties over time. For example, one could write a property stating, "G: if `risk_high` is true, then `autonomy_level ≤ θ`" (where `G` denotes "globally" and `θ` is a predefined threshold). Implementing such a property as a runtime assertion means that the system would actively monitor its state and trigger an error or halt execution if the specified invariant is violated. This transforms the governor from a passive decision-maker into an active verifier of its own adherence to high-level policies. While implementing full-scale model checking of a state machine driven by traces is a more advanced and computationally intensive task, the principle is sound . Model checking algorithms can exhaustively explore all possible states of a finite-state machine to verify if it satisfies a given temporal logic formula [10] . Applying this concept here would involve modeling the autonomy decision process as a state machine where inputs are trace attributes and outputs are decisions. A successful model check would mathematically prove that the governor's logic is free from certain classes of errors regarding its adherence to temporal properties. This represents a significant step up from empirical testing with golden vectors, offering a higher degree of confidence in the system's correctness.

The use of a dedicated ALN shard, such as `augdoctor.autonomy.trace.policy.aln`, provides a powerful architectural pattern for managing these invariants at a governance level . Instead of encoding these rules imperatively within application code, they are declared in a structured, human-readable format that can be parsed and enforced by the system itself. This shard could explicitly declare the mandated decision order, the names and types of fields allowed for cross-platform export, and which fields are considered host-local only . By linking this shard to other configuration files like `neuromorph-eco-profile.aln`, the system creates a declarative dependency graph. A change in the eco profile, for instance, could automatically trigger a re-evaluation of the acceptable ranges for `avg_eco_energy_nj` within the trace policy, ensuring that the exported telemetry always aligns with the current operational constraints . This approach decouples the definition of policy from its implementation, making the system more transparent, auditable, and adaptable. Proving in code that no trace field can influence core balances like BRAIN/WAVE/BLOOD/OXYGEN/NANO directly is a crucial part of this governance, reinforcing the separation of telemetry data from direct control mechanisms . This layered defense-in-depth strategy—combining strict schema enforcement, functional test harnesses, explicit logic ordering, and declarative governance—provides a comprehensive framework for achieving the essential goal of cross-language semantic parity.

| Component | Technology Stack | Primary Function |
|---|---|---|
| `autonomy-trace.v1.json` Schema | All Stacks (Rust, JS, ALN) | Defines the mandatory structure, data types, and value constraints for `AutonomyTraceAttributes` . |
| Golden Test Vectors | All Stacks (Rust, JS, ALN) | Pre-defined input data sets for validating functional equivalence of the autonomy decision logic . |
| Semantic Test Harness | All Stacks (Rust, JS, ALN) | A tool to execute golden test vectors and assert bit-identical `autonomy_level` and `shot_level` decisions . |
| Decision Logic Implementation | Rust, JavaScript | Code that enforces the strict priority order: `identity → consent → lifeforce → eco → SCALE/turns` . |
| Runtime Assertions | Rust, JavaScript | Embedded checks that validate temporal-logic properties of the governor's decisions in real-time . |
| Governance Shard (`augdoctor.autonomy.trace.policy.aln`) | ALN | Declares which trace fields are globally exportable vs. host-local only, creating a verifiable contract . |

This systematic approach ensures that the foundational layer of the system is robust and predictable, providing a stable platform upon which adaptive and learning components can be safely built.

# Governance Architecture for Hard Constraints and Observational Equivalence

The governance architecture for the autonomous bio-feedback system is predicated on a philosophy of safety-first design, where the preservation of the host's well-being is paramount and takes absolute precedence over optimizing for task performance or reward. This is operationalized through a sophisticated system of hard constraints, a veto-first "inner-ledger," and a nuanced verification process focused on "observational equivalence" rather than direct behavioral cloning. The design moves beyond simple access control to create a dynamic, self-enforcing environment where policy is defined declaratively and automatically adapts to changing environmental conditions. This section details the components of this architecture, explaining how hard constraints are implemented as non-negotiable vetoes, how governance is managed via ALN shards, and how a learned policy's success is measured against the gold standard of the `CompanionAutonomyGovernor`.

At the heart of the governance model is the concept of the inner-ledger as a "veto-first oracle" . This metaphor frames the system's core constraints not as soft guidelines that can be negotiated or bent for the sake of efficiency, but as immutable laws that govern all

decision-making processes. The hard constraints explicitly mentioned are lifeforce bands, eco budgets, identity-drift limits, daily turn caps, and a strict self-only scope . Each of these represents a critical boundary that, if crossed, could have detrimental effects on the host's biological or cognitive state. For example, the `lifeforce` constraint ensures the host's physiological resources are not depleted below a safe minimum, while the `eco_budget` constraint manages computational resource consumption to prevent excessive heat generation or energy drain, particularly relevant in BCI contexts. The `identity_drift` limit is a neurorights guardrail designed to prevent the system from making evolutionary changes that could alter the host's sense of self in undesirable ways. The daily turn cap acts as a throttle to prevent runaway optimization loops, and the self-only scope ensures that the agent's actions are confined to its own domain, preventing unintended interactions with external entities. Any policy, whether hand-coded or learned, must prove that it will never propose an action that violates any of these constraints under any circumstances . This makes the inner-ledger the ultimate arbiter of admissible actions.

To manage these constraints effectively, the system leverages a declarative governance model centered on ALN (Agent Language Network) shards . A dedicated shard, tentatively named `augdoctor.autonomy.trace.policy.aln`, would be created to formally declare which fields within the `AutonomyTraceAttributes` object are permitted for cross-platform export and which are strictly host-local . This is a powerful mechanism for separating global telemetry from private, local state. Fields like `stress` or `reward` might be deemed exportable for analysis, while sensitive internal variables like `internal_state_vector` would be marked as host-local only. This shard acts as a verifiable contract that can be programmatically checked by any component in the system, ensuring that data sharing respects privacy and security boundaries. The architecture's dynamism is further enhanced by linking this policy shard to other configuration files, such as `neuromorph-eco-profile.aln` and `assistant-autonomy-profile.aln` . This creates a dependency chain where a change in an environmental profile—such as tightening the available eco bands due to a power-saving mode—is automatically reflected in the policy shard. Consequently, the acceptable range for `avg_eco_energy_nj` in any exported trace would be dynamically adjusted, preventing the system from generating telemetry that is inconsistent with the current operational reality. This automated propagation of constraints significantly reduces the risk of misconfiguration and ensures that the system's policies remain aligned with its physical and environmental context.

The verification of a newly developed learning module introduces a subtle but critical distinction between direct comparison and observational equivalence . The `CompanionAutonomyGovernor` is described as a hand-coded baseline, presumably

representing a carefully crafted set of rules intended to be helpful and effective . However, the research goal stipulates that performance versus this governor is a secondary lens for evaluation . The primary success metric is the learning module's ability to adhere to the hard constraints. A learned policy is considered admissible only if it is observationally equivalent to the governor, meaning it produces the same results *under the same constraints*, without weakening the guards or producing extra, unaccounted-for evolutionary effects . This is a much stronger form of verification than simple behavioral cloning. It requires that the learning process discovers a policy that operates entirely within the same "legal" space defined by the inner-ledger. To facilitate this, the system logs any divergences between the governor's decisions and the learned policy's proposals . These logged cases are treated as valuable research data, allowing engineers to analyze why the learned policy suggested a different course of action. Such analysis could reveal opportunities for improvement in the governor's heuristics or expose limitations in the reward function used for training the learning module. This process fosters a symbiotic relationship between the static, rule-based governor and the dynamic, learned policy, where both are held to the same high standard of constraint adherence.

Furthermore, the governance framework specifies that autonomy traces are to be treated as non-financial, non-identity telemetry, usable exclusively for host-local assistive tuning . This declaration is reinforced by a proof in the code that no field within an autonomy trace can directly influence the host's core balances (BRAIN/WAVE/BLOOD/OXYGEN/NANO) . This is a crucial design choice to prevent the feedback loop from becoming self-referential or manipulative. The traces provide information about the state of the system, which can then be used to adjust the governor's parameters locally, but they cannot directly command a change to the underlying state that produced them. This separation ensures that the system's evolution is guided by an external, observable signal (the trace) rather than being subject to internal feedback loops that could lead to instability or unforeseen consequences. This entire governance architecture, from the veto-first oracle to the declarative ALN shards and the focus on observational equivalence, is designed to build a system that is not only intelligent and adaptive but also provably safe, transparent, and aligned with the long-term interests of the host.

# Empirical Signal Calibration Across Operational Environments

While establishing semantic parity provides a consistent decision-making framework, the quality and relevance of the inputs driving those decisions are paramount. The research

roadmap mandates the empirical calibration of signal ranges for `AutonomyTraceAttributes` across three distinct operational environments: chat-only, neuromorph.softwareonly, and bci.hci.eeg . This calibration is not a one-time setup but an ongoing process of mapping abstract, normalized signals to concrete, host-local phenomena. It involves translating high-level concepts like `lifeforce_scalar` and `eco_energy_nj` into tangible metrics that reflect the host's actual physiological state and the system's computational load. This section outlines the methodology for this crucial calibration process, detailing how abstract signals are grounded in reality through correlation with raw biosignals and hardware counters.

The first major calibration task is to map the `lifeforce_scalar` attribute to existing `LifeforceBandSeries` and `HostEnvelope` fields . The goal is to interpret NeuralRope traces back into meaningful bioscale envelopes without ever exposing the raw, unprocessed brain, blood, or oxygen data directly to the autonomy logic . This requires a two-step process. First, during an initial setup or periodic recalibration phase, the system would collect synchronized data streams: the `lifeforce_scalar` output from the autonomy governor and the corresponding raw biosignals from sensors (e.g., EEG, fNIRS, PPG). Using statistical and machine learning techniques, a mapping function can be derived to correlate fluctuations in the scalar with specific patterns or amplitudes in the raw data. For example, a sustained drop in `lifeforce_scalar` might be correlated with increased theta wave activity in the EEG, indicating mental fatigue. Second, this mapping function, along with the resulting `LifeforceBandSeries` (representing low, medium, high bands), becomes part of the host's personal profile. The `HostEnvelope` fields would then encapsulate these calibrated bands, providing the autonomy logic with clear, pre-defined thresholds (e.g., `min_lifeforce = 0.3`). This ensures that the abstract concept of "lifeforce" is consistently interpreted based on the individual host's unique physiology, moving it from a purely theoretical construct to a measurable and actionable parameter.

The second critical calibration task is to give concrete meaning to the `eco_energy_nj` unit . The roadmap specifies that "1.0" should represent a concrete, host-local high-band value, not an arbitrary constant. This requires a tight integration between the software-level energy accounting and the physical hardware's power consumption. The process begins by instrumenting the system to track `FLOP` (Floating Point Operations Per Second) counters and other relevant computational metrics. Controlled experiments would then be conducted where a fixed computational workload is executed, and the corresponding increase in energy consumption is measured using hardware sensors (e.g., power meters, thermal cameras). This establishes a calibration factor that translates FLOPs into joules of energy consumed. The `eco_energy_nj` value recorded in an `AutonomyTraceAttributes` object would then be calculated by multiplying the

number of FLOPs executed by this calibration factor. The `neuromorph-eco-profile.aln` file would define the total energy budget available to the host . During operation, the system would maintain a running tally of `eco_energy_nj` spent, comparing it against the `eco cap` derived from the profile. This end-to-end calibration links the abstract notion of an "eco-cost" to a physically meaningful quantity, allowing the system to make informed decisions about computational trade-offs in a way that respects the host's physical and environmental constraints.

Defining canonical ranges for all `AutonomyTraceAttributes` fields—such as `stress`, `fatigue`, `reward`, `safety`, and `risk`—is another key aspect of the calibration plan . These ranges are necessary for normalizing inputs to the learning models and for setting clamping bounds during processing. The process for determining these ranges will vary depending on the signal's nature and the operational environment.

For the **chat-only environment**, many of these signals may be proxies derived from text analysis and user interaction patterns.

- `Stress` and `Fatigue`: Could be estimated from typing speed, error rates, sentence complexity, and message length.
- `Reward`: Could be derived from task completion success rates or positive sentiment in the conversation history.
- `Risk`: Could be a function of response latency and the presence of certain keywords or phrases indicating uncertainty or danger.

For the **neuromorph.softwareonly environment**, signals would be primarily generated by the internal simulation of the host's neural model.

- `Stress` and `Fatigue`: Could be modeled based on the computational load of the simulation, the novelty of encountered problems, or the entropy of internal state transitions.
- `Reward`: Would be tied to the accuracy of the model's predictions and the coherence of its simulated thoughts.
- `Risk`: Could be calculated based on the divergence of the simulation from a known baseline or the probability of entering an unstable state.

For the **bci.hci.eeg environment**, the signals would be directly derived from neurophysiological measurements.

- `Stress` and `Fatigue`: Could be quantified from EEG features like alpha/beta wave ratios, heart rate variability (HRV) extracted from ECG or PPG, or galvanic skin response (GSR).

- `Reward`: Might be inferred from patterns associated with positive affect, such as increases in gamma-band power or P300 event-related potentials.
- `Risk`: Could be detected from indicators of cognitive overload or conflict, such as elevated frontal midline theta (FMT) activity.

The table below summarizes the proposed calibration targets for each signal across the three environments.

| Signal Attribute | Description | Calibration Target / Methodology (Chat-Only) | Calibration Target / Methodology (Software Simulation) | Calibration Target / Methodology (BCI/HCI/EEG) |
|---|---|---|---|---|
| stress | Subjective or physiological strain. | Text-based proxies (typing speed, error rate, sentiment analysis). | Simulation-based proxies (computational load, state entropy, prediction error). | Direct measurement from EEG (alpha/beta ratio), HRV, GSR. |
| fatigue | Mental exhaustion or drowsiness. | Text-based proxies (sentence length, repetition, pauses). | Simulation-based proxies (long-term state degradation, decreased prediction accuracy). | Direct measurement from EEG (theta power, alpha asymmetry). |
| eco_nj | Computational energy cost in nanojoules. | Estimated from CPU/GPU usage and network I/O during text rendering. | Tied to FLOP counters and memory access patterns of the simulation. | Tied to FLOP counters and power draw of connected hardware. |
| reward | Positive outcome or reinforcement. | Task success, positive user feedback, completion of sub-tasks. | Accuracy of model predictions, coherence of simulated thought streams. | User-reported comfort sliders, P300 ERP amplitude. |
| safety | Perceived or actual safety of the action. | Heuristic rules based on content moderation APIs and context. | Simulation of action consequences within the model's world. | Real-time monitoring of BCI signal stability [28]. |
| lifeforce | Host's physiological/ cognitive resource level. | Proxy for user-reported energy levels or battery percentage. | Internal state variable representing simulation health. | Mapped from raw biosignals (BRAIN/BLOOD/OXYGEN) via a calibrated function . |
| risk | Probability of negative outcome or harm. | Keyword analysis, sentiment of responses, API error rates. | Divergence from stable baselines, probability of state collapse. | Indicators of cognitive conflict or overload (e.g., FMT theta) [21]. |

By systematically performing these calibrations, the system grounds its abstract decision-making process in a rich tapestry of empirical data. This ensures that the autonomy logic is not operating on sterile, meaningless numbers but is instead responsive to the host's true state and the tangible costs of its operations, enabling a truly adaptive and context-aware experience.

# Designing Neurorights-Compliant Assist Ramps for Bio-Feedback Coupling

The integration of bio-feedback into the autonomy system, forming a closed-loop control mechanism, is governed by a stringent requirement for smoothness and neurorights compliance. The central concept for this integration is the "assist ramp"—a function that maps a vector of bio-feedback signals (e.g., stress, fatigue, eco usage) to a suggested `autonomy_level`. The user's directive prioritizes the mathematical and ethical properties of these ramps above their immediate predictive performance. The core requirements are that the ramps must be smooth, exhibiting no discontinuities, and must operate entirely within the defined neurorights envelopes, which prohibit sharp spikes in identity-drift or violations of pain/psychrisk budgets . This approach treats the human operator as a delicate system requiring gradual, predictable adaptation rather than sudden interventions. The design of these ramps is thus a synthesis of signal processing, control theory, and human-computer interaction principles.

The most profound insight guiding this design is the mandate for smoothness, which is explicitly linked to the mathematical concept of Lipschitz continuity . A function $f(x)$ is said to be Lipschitz continuous if there exists a constant $L$ (the Lipschitz constant) such that for any two points $x_1$ and $x_2$, the inequality $|f(x_1) - f(x_2)| \leq L|x_1 - x_2|$ holds. In the context of assist ramps, this property guarantees that the rate of change of the suggested autonomy level is fundamentally bounded. A discontinuity in the ramp would correspond to an infinite Lipschitz constant, which would manifest as an instantaneous, jarring jump in the system's autonomy level. Such a jump could induce cognitive dissonance, stress, or a feeling of loss of control in the user. By restricting the assist ramps to be Lipschitz-bounded functions, such as sigmoidal curves or carefully constructed piecewise-linear functions, the system ensures that all adjustments to autonomy are gradual and predictable . This mathematical guarantee provides a strong foundation for building trust and ensuring a seamless user experience. The choice of the Lipschitz constant $L$ itself becomes a critical hyperparameter, representing the maximum permissible rate of assistance. This value could be informed by literature on neurofeedback, which suggests that gradual transitions are key to effective neural adaptation and skill acquisition [1] [2], or by studies on the plug-and-play stability required for reliable BCI systems [28].

Designing and validating these assist ramps involves a rigorous, multi-stage process. The first stage is the design phase, where candidate functions are formulated. These functions take the normalized bio-feedback signals as input and produce a continuous suggestion for the `autonomy_level`. For example, a simple ramp might increase autonomy slightly as `stress` rises, but only if `fatigue` is low and `eco_nj` is well within budget. The

second stage is the validation phase, which focuses on verifying two key properties: smoothness and budget adherence. Smoothness can be validated analytically by inspecting the derivative of the ramp function to ensure it never exceeds a predefined maximum value (the inverse of the desired Lipschitz constant). Budget adherence is validated by simulating the ramp's operation over extended periods of trace data. The system would track cumulative `identity_drift` and `pain/psychrisk` exposure, ensuring they remain within their respective `neurorights` envelopes . Any simulation run that leads to a budget violation would invalidate that particular ramp design, signaling that its sensitivity to certain inputs is too high. This iterative process of design, simulation, and validation continues until a family of safe and smooth assist ramps is established.

Once a set of validated, neurorights-compliant assist ramps is in place, the third stage is optimization. The user's directive is clear: optimization for high reward and low stress must occur *within* the safe boundaries defined by the validated ramps . This sequential approach—establishing safety first, then optimizing performance—is a hallmark of robust engineering. It avoids the common pitfall of designing a highly performant but dangerously volatile system. The optimization process involves experimenting with different combinations of bio-feedback predictors fed into the ramp functions. For instance, one might compare a ramp that uses only `stress` as a predictor against one that combines `stress` and `HRV proxy`. The performance of each combination is evaluated using historical `NeuralRope` trace data. The goal is to find the predictor set that, when passed through the validated smooth ramp, yields the highest average `reward` while simultaneously minimizing average `stress` and `eco_nj` consumption . This search is conducted entirely within the feasible region defined by the safety constraints, ensuring that no matter how the predictors are tuned, the system's behavior remains predictable and respectful of the host's well-being.

Finally, the system incorporates a mechanism for continuous improvement through active learning, using trace data to identify "negative exemplars" . The `NeuralRope` captures a rich history of system states and outcomes. Segments of this rope where the autonomy level was high but the resulting `reward` was low, or the `risk` score was unexpectedly high, are flagged as negative examples. These instances represent failures of the current policy or ramp functions. They are invaluable for refining the system. For example, if the data shows that the system frequently grants high autonomy to users experiencing mild confusion (high `error_rate` from Reality.os), leading to poor task performance, this pattern can be used to tighten the `risk` or `safety` thresholds in the governor. Similarly, if a particular combination of predictors in an assist ramp consistently fails to respond appropriately to a specific biomarker pattern, the ramp's mathematical form or its underlying weights can be adjusted. This process of mining traces for failure modes and

using them to iteratively refine the governor and assist ramps creates a powerful, self-correcting feedback loop that enhances the system's intelligence and safety over time.

# Autonomous Policy Learning Under Constrained Optimization

The development of an autonomous policy learning component represents the apex of the system's adaptive capabilities. However, the research framework places this objective firmly within a tightly constrained environment, where the primary success metric is not raw performance but strict adherence to the hard constraints defined by the inner-ledger. The learning process is framed as a problem of constrained optimization, where the objective is to learn a policy $\pi(a|s)$—the probability of taking an action `a` (an `autonomy_level` and `shot_level_label`) given a state `s` (the `biomarkers` from the trace)—that maximizes cumulative reward while respecting non-negotiable bounds on risk, lifeforce, and eco usage . This approach, often referred to as safe reinforcement learning, ensures that the system's exploration of the policy space does not come at the expense of the host's safety or well-being.

The formulation of the learning problem is critical. Each `AutonomyTraceAttributes` JSON object is treated as an off-policy transition, containing a complete snapshot of a decision point . The state `s` is defined by the biomarker readings (stress, fatigue, etc.), the action `a` is the chosen autonomy level and shot level, the reward `r` is the observed `avg_reward`, and the primary constraints are derived from the `constraint` field, specifically the `highest_risk_score` and `worst_lifeforce_scalar` . The learning algorithm, whether a shallow Reinforcement Learning (RL) agent or a Constrained Bandit solver, must be explicitly designed to incorporate these constraints into its update rules or action selection process . For instance, in a Constrained Multi-Armed Bandit formulation, the agent would not simply choose the arm (action) with the highest expected reward but would select from a subset of arms that are known to satisfy the constraints on risk and lifeforce. In a more general RL setting, techniques such as Lagrangian relaxation or primal-dual methods can be used to penalize constraint violations in the objective function, effectively teaching the policy to avoid dangerous regions of the state-action space.

The implementation of this learning module in either Rust or JavaScript must be architected to interface seamlessly with the core governance structures . The agent's knowledge of the constraint boundaries—the `max_risk`, `min_lifeforce`, and `eco`

cap—must be sourced directly from the `neuromorph-eco-profile.aln` and other relevant governance shards . This ensures that the learning process is always working with the latest, officially sanctioned limits, and prevents the policy from drifting out of alignment with the host's configured safety settings. The learning algorithm would operate on the normalized signal ranges that were established during the empirical calibration phase, ensuring that its inputs are consistent and meaningful. The output of the learning process is not just a new policy, but a new hypothesis about how to best allocate autonomy under pressure. This hypothesis must be subjected to the same rigorous verification process as any other component of the system.

Verification of the learned policy is a two-pronged process centered on constraint adherence and observational equivalence. As previously discussed, the primary verification step is to confirm that the learned policy never proposes an action that would be rejected by the lifeforce/eco guards . This is a binary pass/fail condition. If the policy ever suggests a state where `lifeforce < min_lifeforce` or `eco_nj > eco_cap`, the policy is immediately discarded. This is the "veto-first" principle in action. The secondary verification step involves comparing the learned policy's decisions against those of the hand-coded `CompanionAutonomyGovernor` . This comparison is not about finding the policy that performs best in terms of raw reward, but about ensuring that the learned policy is "observationally equivalent." This means that for any given state `s`, the learned policy's distribution over actions `a` should be indistinguishable from the governor's deterministic choice, provided both are operating within the same constraint boundaries. Any divergence is logged as a research case for analysis . These cases are particularly interesting when the learned policy chooses a different action than the governor. If the governor's action leads to a high-risk outcome, the learned policy's choice might represent an improvement. Conversely, if the governor's action was correct, the divergence might indicate a flaw in the learned policy's reward function or its understanding of the constraints. This comparative analysis helps to build a corpus of insights that can be used to further refine both the governor and the learning algorithm.

The ultimate goal of this constrained learning process is to augment, not replace, the `CompanionAutonomyGovernor`. The governor provides a baseline of "helpfulness" and a proven, rule-based fallback . The learned policy, trained on vast amounts of trace data, has the potential to discover more nuanced and efficient strategies for allocating autonomy. However, its value is only realized if it can demonstrate that it adheres to the same uncompromising standards of safety and constraint adherence as the governor. By framing the learning problem this way, the research ensures that the pursuit of intelligence is always subordinate to the preservation of the host's integrity, creating a hybrid system where the strengths of both rule-based reasoning and data-driven learning are harnessed in a safe and synergistic manner.

# Integrated Trace Analysis for System Verification and Refinement

The `NeuralRope` and the `AutonomyTraceAttributes` objects it contains serve as the system's collective memory, providing an invaluable, high-fidelity record of its operation over time. This trace data is not merely a diagnostic tool but the central substrate for the entire research and refinement cycle. It is the evidence base used to verify semantic parity, validate assist ramps, compare policies, and identify opportunities for improvement. The systematic analysis of this data is therefore a critical, recurring theme that ties together all aspects of the research framework. This section details how trace data is structured, captured, and analyzed to drive the continuous verification and refinement of the autonomous bio-feedback system.

The structure of an `AutonomyTraceAttributes` object is designed to capture the essential elements of a decision-making moment, formatted as a JSON object suitable for logging and analysis . A typical trace entry includes several key fields:

- **State Variables:** These represent the host's physiological and cognitive state at the moment of decision. Key fields include `stress`, `fatigue`, `lifeforce_scalar`, `eco_nj`, `safety`, and `risk` . These are the normalized, calibrated signals that feed into the decision logic.
- **Action Taken:** This field records the decision made by the system. It is composed of the `autonomy_level` granted and the `shot_level_label` selected . This is the output of the governor or learning policy.
- **Outcome Metrics:** These fields describe the result of the action. The primary metric is `reward`, representing the utility gained from the action. Other fields like `avg_reward` provide a rolling average for trend analysis .
- **Constraint Violation Flags:** This is a crucial metadata field, `constraint`, which captures the severity of the situation. It is typically defined by the `highest_risk_score` and `worst_lifeforce_scalar` observed during the transition, providing a clear signal of any constraint breaches .
- **Source Context:** Information identifying the source of the trace, such as the `host_id` and the `operating_environment` (e.g., 'chat-only', 'bci.hci.eeg'), is also included to allow for environment-specific analysis.

The recorder component is responsible for ingesting and persisting these trace objects, extending its capabilities to include additional bio-feedback channels as needed . These can include proxies for HRV, subjective comfort ratings from sliders, and even error rates

from an external system like `Reality.os`. This enriched data stream provides a more holistic view of the host-system interaction, enabling more sophisticated analyses.

The primary application of trace analysis is in the verification of cross-language semantic parity. The "semantic test harness" relies on replaying sequences of trace data through the governors in Rust and JavaScript. By asserting bit-identical decisions on a large corpus of diverse traces, the system can empirically validate that the semantic invariant holds true across a wide range of scenarios, far exceeding the coverage of manually crafted golden test vectors. This continuous verification process acts as a safety net, catching regressions in the decision logic whenever the codebases are updated.

A second major application is in the validation and refinement of assist ramps. After a ramp is designed, it is applied to a segment of trace data to simulate its performance. The resulting output is compared against the original `autonomy_level` decisions to see if the ramp produces a similar, yet smoother, trajectory. The analysis focuses on identifying any discontinuities or budget violations that the ramp might introduce. Furthermore, trace data is mined for "negative exemplars"—segments where the system's previous behavior led to poor outcomes. For example, an analyst might filter the traces for all instances where `autonomy_level` was high but `reward` was subsequently low. These segments represent failure modes. By examining the state variables (`stress`, `fatigue`, etc.) in these negative exemplars, engineers can diagnose weaknesses in the current governor thresholds or assist ramp sensitivities. If the data consistently shows that high autonomy is ineffective when `fatigue` is above a certain threshold, the governor's logic can be updated to prevent granting high autonomy in that regime. This data-driven approach to refinement ensures that improvements are targeted and based on real-world evidence of system performance.

Finally, trace analysis is the bedrock of the constrained policy learning and verification process. The dataset of traces forms the basis for training the learning module. The analysis of the traces' `constraint` fields is used to pre-filter the data, creating a clean dataset of "safe" transitions on which the policy can be trained. During and after training, the learned policy's proposed actions are compared against the `action_taken` field in the traces. The logged divergences between the learned policy and the `CompanionAutonomyGovernor` are a primary source of research material. By clustering these divergent cases based on their state variables, researchers can understand the specific conditions under which the learned policy explores different strategies. Are these exploratory actions safer? More efficient? Or do they represent a genuine failure? Answering these questions requires deep dive analysis of the surrounding trace segments, looking at the preceding and succeeding states to understand the causal chain of events. This entire workflow—from data capture to

structured analysis—transforms the `NeuralRope` from a simple log into a dynamic, living document that drives the evolution of the system, ensuring that every change is grounded in empirical evidence and rigorously verified against the core principles of safety and semantic consistency.

---

### Reference

1. An Active Inference perspective on Neurofeedback Training https://arxiv.org/pdf/2505.03308?

2. Optimal Inter-Session Intervals in Neurofeedback Training https://www.mdpi.com/2076-3417/16/1/142

3. A Neuroadaptive Blueprint for Non-Invasive Vision ... https://www.researchgate.net/publication/390345115_Through_the_Ear_We_See_A_Neuroadaptive_Blueprint_for_Non-Invasive_Vision_Restoration_via_Auditory_Interfaces_Proposed_Table_of_Contents_Part_I_Foundations_of_Sensory_Rerouting

4. 目录(Contents) https://ieeexplore.ieee.org/iel7/6621666/6639389/06641025.pdf

5. Keum Shik Hong https://www.sciencedirect.com/author/7402515710/keum-shik-hong

6. IEEE Project Domains 2020 - 2021: We Can Help You To ... https://www.scribd.com/document/542296016/Ece-Project-Topics

7. IEEE Project Titles 2020 | PDF | Android (Operating System) https://www.scribd.com/document/502413902/IEEE-Project-Titles-2020

8. Computer Science https://arxiv.org/list/cs/new

9. Soft Computing and Signal Processing https://link.springer.com/content/pdf/10.1007/978-981-97-9926-8.pdf

10. Fundamental Results for the Verification of Observational ... https://www.researchgate.net/publication/221248621_Fundamental_Results_for_the_Verification_of_Observational_Equivalence_A_Survey

11. High-Entropy Oxide Memristors for Neuromorphic Computing https://pmc.ncbi.nlm.nih.gov/articles/PMC12378895/

12. Roadmap to Neuromorphic Computing with Emerging ... https://arxiv.org/html/2407.02353v1

13. Neuromorphic Hardware for Artificial Sensory Systems https://link.springer.com/article/10.1007/s11664-025-11778-x

14. Van der Waals integration of 2D materials for advanced ... https://iopscience.iop.org/article/10.1088/2634-4386/ae294e

15. 30 years of AI for electrocatalysis: Where we are and what's ... https://www.sciencedirect.com/science/article/pii/S2667141725001454

16. 2022 roadmap on neuromorphic computing and engineering https://hal.science/hal-03872100v2/file/Christensen_2022.pdf

17. Autonomous driving controllers with neuromorphic spiking ... https://www.frontiersin.org/journals/neurorobotics/articles/10.3389/fnbot.2023.1234962/full

18. NOVEMBER 2019 https://www.ieee.org/ns/periodicals/NxtBooks/SP/PDF/SP_Nov2019.pdf

19. The h-BCI system interfaces. (A) WhatsApp ... https://www.researchgate.net/figure/The-h-BCI-system-interfaces-A-WhatsApp-interface-contact-up-and-contact-down_fig1_361088561

20. 2024-31-01 Database Search TLR Reliance https://hal.science/hal-04637901v1/file/2024-31-01%20Database%20Search%20TLR%20Reliance.pdf

21. Advances in Brain Inspired Cognitive Systems - Springer Link https://link.springer.com/content/pdf/10.1007/978-981-96-2885-8.pdf

22. 机器学习2025_5_23 http://www.arxivdaily.com/thread/67749

23. Assistive Tech Companies Worldwide and in India With ... https://www.linkedin.com/pulse/assistive-tech-companies-worldwide-india-gaps-m5n0e

24. Annual Review of Cybertherapy and Telemedicine https://interactivemediainstitute.com/wp-content/uploads/2019/04/ARCTT-7.pdf

25. Chinese Control and Decision Conference (CCDC 2021) https://ieeexplore.ieee.org/iel7/9601162/9601344/09602651.pdf

26. Deep Learning Adapted to Differential Neural Networks ... https://www.researchgate.net/publication/350168065_Deep_learning_adapted_to_differential_neural_networks_used_as_pattern_classification_of_electrophysiological_signals

27. Reinforcement Learning: An Introduction | Guide books https://dl.acm.org/doi/book/10.5555/3312046

28. Plug-and-Play Stability for Intracortical Brain-Computer ... https://www.researchgate.net/publication/375793747_Plug-and-Play_Stability_for_Intracortical_Brain-Computer_Interfaces_A_One-Year_Demonstration_of_Seamless_Brain-to-Text_Communication

29. NeuroChat: A Neuroadaptive AI Chatbot for Customizing ... https://dl.acm.org/doi/10.1145/3719160.3736623

30. Gait Recognition for Lower Limb Exoskeletons Based on ... https://pmc.ncbi.nlm.nih.gov/articles/PMC8976668/