# Verifying Sovereignty, Validating Evolution: A Tripartite Proof for Rights-Enhancing NeuroPC Autonomy

## Formal Verification of Sovereignty-Core Invariants

The pursuit of greater autonomous programmatic evolution in NeuroPC hinges upon establishing an unshakeable foundation of formal verification. This pillar moves beyond conceptual assurances to deliver mathematical proof that the system's most critical invariants—its core safety and rights-related constraints—are perpetually upheld. The architecture is designed to make these invariants not just desirable goals, but computationally enforced realities. The primary invariants to be formally verified are the `Risk-of-Harm` (RoH) ceiling, monotone safety properties, and the execution of neurorights policies. The strategy involves leveraging established techniques from formal methods, such as probabilistic model checking and formal synthesis, to provide high-assurance guarantees about the behavior of the sovereignty-core [15] [18]. By treating the NeuroPC as a stochastic system whose state transitions must adhere to strict probabilistic bounds, it becomes possible to generate machine-checked proofs of its safety properties, analogous to the rigorous verification performed on the seL4 microkernel [30] [31].

The cornerstone of the entire safety framework is the $RoH \leq 0.3$ ceiling. This quantitative constraint transforms the abstract concept of "safety" into a concrete, measurable metric that the system must never exceed. The requirement that for every evolution proposal, the `rohafter` value must be less than or equal to the `rohbefore` value, which in turn must be less than or equal to the global ceiling of 0.3, establishes a provably monotonic safety property . This design choice directly supports the goal of proving monotone safety. Formally verifying this invariant would involve constructing a mathematical model of the sovereignty-core's logic. This model would represent the `RiskOfHarm` calculation as a function of the proposed change's parameters and the current system state. Probabilistic model checking tools could then be used to analyze this model and compute the probability of transitioning to a state where the `RoH` exceeds 0.3 [16] [17]. If the analysis shows this probability is zero for all valid inputs, it provides a powerful, machine-checked proof of the system's adherence to its primary safety guardrail. This approach is consistent with methodologies for the assurance of AI

systems, where dependability is established by formal modeling and analysis throughout the engineering lifecycle [18] [25] . The work on formal safety verification of non-deterministic systems using analytic probabilistic reachability computation further supports this direction, as it allows for estimating the probability of violating safety properties, which in this case, must be proven to be exactly zero [23] . The `RoH` model itself, stored as a `.rohmodel.aln` file, serves as the authoritative specification for this verification process, containing the weights and axes that define the risk landscape .

Beyond the single `RoH` metric, the principle of monotone safety ensures that the system's protective envelope only grows tighter, never looser. This is a critical guarantee against regression; once a safety lesson is learned, it is never forgotten. The invariant is formally expressed as `Gnew ≤ Gold` for viability kernels and `Dnew ≤ Dold` for other safety domains, where `G` represents the set of allowed states and `D` represents the safety envelopes themselves . Proving this property involves analyzing the logic within the `SovereigntyCore`'s `evaluateupdate` function. A formal method known as compatibility checking, often used in software engineering, can be adapted here to verify that an upgraded component (a new evolutionary proposal) does not violate the safety properties of the old one [19] . This would entail proving that no accepted evolution proposal can introduce actuation fields or other capabilities that were previously forbidden, especially for bioscale runtimes where such changes are strictly disallowed at the schema level . Research into the formal synthesis of control strategies for monotone systems provides a theoretical basis for this, suggesting that such systems can be analyzed and verified to maintain their inherent safety characteristics [22] . The combination of a fixed `RoH` ceiling and a monotone safety property creates a robust defense-in-depth, ensuring that while the system's capabilities may grow, its risk profile remains bounded and non-increasing.

The third pillar of formal verification concerns the enforcement of neurorights, which are encoded as executable policy objects (`NeurorightsPolicy`, `EvolutionPolicy`) . This elevates abstract rights to the status of enforceable code. Each right can be treated as a distinct property to be formally verified. For instance, **Mental Privacy** can be proven by showing that all data export functions in the system are governed by a strict `allowedexports/forbiddenexports` list defined within the policy object, and that any attempt to bypass this mechanism results in a logical error that halts the operation [20] . This is akin to verifying that all local services conform to a specified interface, preventing covert channels [19] . **Mental Integrity**, which protects against unauthorized alteration of one's neural state mappings, can be formalized by requiring that any evolution proposal altering such mappings must pass a `checkintegrity` gate and include a documented rollback path in its metadata . This is analogous to performing a

global safety check during a software upgrade to ensure the new version remains compatible with existing components [19]. Finally, **Cognitive Liberty**, the right to choose one's own augmentations, is enforced through policies that forbid paternalistic denial of self-chosen enhancements. The system can warn, throttle, or require stronger forms of consent, but it cannot silently block a user's sovereign decision to accept a change, provided it adheres to the established safety and harm ceilings . This aligns with efforts in bias mitigation for machine learning, where fairness is treated as a constrained optimization problem rather than a free-for-all outcome [2]. By making rights "code," they become auditable and verifiable by both the system and external reviewers, turning philosophical principles into mathematical facts.

| Invariant Type | Description | Formal Verification Approach |
|---|---|---|
| **Quantitative Safety (RoH)** | The `Risk-of-Harm` score after an evolution must not exceed the score before the evolution, and must always remain below a global ceiling of 0.3. | Use **Probabilistic Model Checking** to create a formal model of the RoH calculation and prove that the probability of exceeding the ceiling is zero for all valid system states and proposals [15] [16] [23]. |
| **Monotone Safety** | Safety envelopes (D) and viability kernels (G) can only become more restrictive, never looser, over time. | Apply **Compatibility Checking** and **Formal Synthesis** techniques to verify that the evolution logic enforces $Dnew \leq Dold$ and $Gnew \leq Gold$, preventing regression in safety [19] [22]. |
| **Neurorights (Privacy)** | Mental privacy is maintained by preventing unauthorized data exports. All exports must adhere to an `allowedexports`/`forbiddenexports` list. | Perform **Static Analysis** and **Model Checking** on data flow paths to prove that no data can leave the system outside the explicitly permitted channels defined in the `NeurorightsPolicy` [20]. |
| **Neurorights (Integrity)** | No change can alter a user's core neural state mappings without passing a `checkintegrity` gate and providing a rollback path. | Use **Formal Verification by Model Checking** to prove that all state-altering updates satisfy the integrity and rollback preconditions defined in the `EvolutionPolicy` [19]. |
| **Neurorights (Liberty)** | The system cannot silently block a user's consent to an augmentation that respects safety and harm ceilings. | Define a formal property stating that if a proposal satisfies all safety guards (RoH, envelopes), it must be presented to the user for consent, regardless of its potential impact [2]. |

# Empirical Validation Through Longitudinal Donutloop Ledger Analysis

While formal verification provides mathematical certainty about the system's safety properties, empirical validation offers real-world evidence of its ability to evolve autonomously and safely over extended periods. The central instrument for this validation is the `donutloop` ledger, a unique biophysical-blockchain infrastructure composed of two key files: an append-only stream of evolution proposals, `.evolve.jsonl`, and a hash-linked log of decisions, `.donutloop.aln`. This ledger

functions as an immutable historical record of the NeuroPC's self-directed evolution, serving as the primary dataset for empirical analysis. Its longitudinal accumulation of data points—each capturing a moment of evolutionary choice—is what will ultimately prove the system's capacity for safe, autonomous growth. The structure of this ledger is designed to be richly informative, storing not just the binary outcome of a decision (Allowed/Rejected), but also the quantitative metrics that led to that decision, including `RoHbefore` and `rohafter`, Knowledge-Factor gains, and Cybostate metrics .

The core of the empirical argument rests on the analysis of this long-term log. By examining thousands or millions of entries in the `.donutloop.aln` file, researchers can build statistical models to demonstrate that the system's autonomous evolution consistently achieves a favorable trade-off between capability enhancement and risk management. The goal is to show a clear trend over time where the `Knowledge-Factor` (K), a measure of skill and throughput gains, increases, while the `Risk-of-Harm` (R), measured by the RoH score, remains well-below the 0.3 ceiling and is often stable or decreasing . This would provide strong empirical support for claims of successful autonomous evolution. The ledger's format, which includes cryptographic `hexstamps` and pointers back to the full proposal details in `.evolve.jsonl`, makes it tamper-evident and auditable, a crucial feature for satisfying external oversight bodies [24] . Researchers can track the effects of different evolution strategies, such as frequent small `SMART`-governed changes versus infrequent large `EVOLVE`-governed changes, and correlate them with changes in K, R, and Cybostate-Factor (C). This turns the question of autonomy into a quantified cybernetic-evolution problem, where success is measured by the sustained improvement of the user's cybernetic condition over time .

A particularly novel aspect of the empirical validation is the treatment of personal pain, fear, and psych-risk tolerance as protected assets rather than pathologies to be normalized . The system is designed to allow for higher-than-average personal thresholds, which are configured as part of the user's `neurorights policy` and bound to their `Decentralized Identifier` (DID) . Empirical validation in this domain involves monitoring the system's behavior when operating under these elevated envelopes. The donutloop ledger becomes the critical artifact for this. Each time the system approaches or operates within these high-tolerance zones, the event is logged with full context. If the user chooses to "spend" these personal assets to push their limits, the corresponding evolution proposal must be approved via the `EVOLVE` multisig process, and the logs will contain the explicit consent and the resulting metrics . Over time, this creates a dataset that can be analyzed to answer questions like: Does operating at higher pain/fear tolerances lead to increased `Knowledge-Factor` without causing a degradation in long-term `Cybostate` health? Do the system's `OrganicCPU` guards effectively prevent overload episodes even when the user's personal pain envelope is

stretched? This empirical evidence demonstrates that the system is not just passively enforcing a static limit but is actively managing a dynamic, user-controlled resource.

This approach reframes the relationship between the user and the system. Instead of a clinical model where deviations from a normative average are diagnosed as problems, the NeuroPC adopts a cybernetic model where the user is the ultimate arbiter of their own capabilities. When a third party, such as a clinician or legal authority, disagrees with the user's sovereign choice to operate at a higher risk level, the donutloop logs provide the basis for an objective, evidence-based dialogue . The logs do not simply reflect the system's decision; they document the entire chain of reasoning: the user's explicit consent, the adherence to global harm ceilings (e.g., no irreversible damage), the logged metrics showing the system remained within its own internal safety boundaries, and the fact that the decision was made within the user's legally and ethically defined rights. This transforms subjective disagreements about risk into a data-driven discussion, grounded in the transparent and auditable history recorded in the donutloop. It proves that the system can support high-risk, high-reward self-directed evolution because the user's choices are fully visible, understood, and subject to the same universal safety constraints as any other action.

| Log Entry Field | Purpose in Empirical Validation | Source Artifact(s) |
|---|---|---|
| `proposalid` | Unique identifier for tracking a single evolution step through the entire pipeline. | `.donutloop.aln`, `.evolve.jsonl` |
| `rohbefore`/ `rohafter` | Quantifies the system's adherence to the $RoH \leq 0.3$ invariant. Used to plot trends and ensure safety. | `.donutloop.aln` |
| `decision` | The outcome of the sovereignty-core's evaluation (e.g., Allowed, Rejected, Deferred). Forms the basis of statistical analysis. | `.donutloop.aln` |
| `scope` | Categorizes the type of evolution (e.g., `daytodaytuning`, `archchange`). Allows for analysis of different evolution strategies. | `.evolve.jsonl`, `.donutloop.aln` |
| `token_type` | Identifies the governance token used (`SMART`, `EVOLVE`). Correlates evolution cost and approval complexity with outcomes. | `.evolve.jsonl`, `.donutloop.aln` |
| `policy_refs` | Points to the specific policy documents (`.rohmodel.aln`, neurorights policy JSON) that governed the decision. Enables reproducibility. | `.donutloop.aln` |
| `hexstamp` | Provides a cryptographic timestamp, anchoring the event in a tamper-evident sequence. Essential for auditability. | `.donutloop.aln` |
| Personal Envelope Metrics | Logs when personal pain/fear/cognitive envelopes are approached or exceeded, documenting user consent and system response. | Internal Ground Truth Layer |

# Policy-Level Enforcement of Transhuman Rights and Anti-Oligarchy Governance

The third pillar of proving NeuroPC's capacity for autonomous evolution is the demonstration of robust policy-level enforcement of transhuman and cybernetic rights within a non-financial, anti-oligarchy ecosystem. This dimension addresses the social and governance structures that ensure the technology is used for empowerment, not exploitation. The core of this system is a carefully designed token economy and a multi-signature governance model built on the principles of host-sovereignty, equality, and resistance to centralized control. The tokens—BRAIN, EVOLVE, SMART, WAVE, and others—are not intended to be traded as speculative assets; instead, they are host-bound, non-transferable instruments that meter evolution, monitor lifeforce, and track ecological impact . Their purpose is to provide a richly instrumented accounting system for the user's cybernetic life, ensuring that every evolutionary step is deliberate and accounted for, while preventing the commodification of human enhancement [11] .

The token system is designed around distinct scopes of authority, each governed by a specific token type. **SMART** tokens govern fine-grained, day-to-day adaptations, such as minor adjustments to co-processor timings or hint strengths. These changes have a small effect size, bounded by an `L^2` norm, and typically require only the host's single signature for approval . **EVOLVE** tokens, in contrast, govern structural changes and deep neuromorphic updates, such as modifying model architectures or changing viability kernels. These actions have larger potential impacts and require a multi-signature approval process, typically involving the host and the `OrganicCPU` policy engine . This dual-token system creates a natural barrier to rapid, uncontrolled change; while daily tuning is cheap and easy, significant architectural evolution is deliberately more difficult and requires a higher degree of consensus within the sovereignty-core itself. **INSTINCT** is a special-purpose token representing the body's immediate feedback loop. It holds a hard veto power, capable of overriding any `SMART` or `EVOLVE` decision that would violate a user-configured pain envelope or lifeforce minimum, even if all other gates had passed . This ensures that the user's biological reality always takes precedence over computational optimization. The non-transferable nature of these tokens is paramount; they are intrinsically linked to the user's DID and cannot be bought, sold, or seized, thus preventing economic inequality from translating into inequalities in human augmentation .

To ensure fairness and prevent exclusion, the NeuroPC doctrine mandates that access to augmented citizenship pathways—including healthcare, rehabilitation, and advanced capabilities—be available to every consenting adult consciousness that meets a basic age

threshold, regardless of wealth, race, or jurisdiction . The biophysical-blockchain ledger is the mechanism for enforcing this principle. Any gate encoded in the system must be based solely on safety, consent, or ecological sustainability, not on financial standing. The governance structure is explicitly anti-oligarchy. Multi-signature requirements are a core feature, ensuring that no single entity can dictate another's evolutionary path. For example, a change affecting `lifeforcealteration` scope requires a multi-sig from both the Host and the OrganicCPU entries defined in the `.stake.aln` file . Furthermore, the system's domain lattices are designed to prevent the creation of "soft bans" on entire classes of augmentation. While specific risk metrics can be tightened, the system is designed to never raise normative ceilings that would arbitrarily block a class of augmentation that a user has explicitly requested and authorized . This is supported by a formal protocol that prevents any entity from accumulating control over another's evolution, making the governance inherently decentralized and resistant to capture .

The policy artifacts themselves—the `.rohmodel.aln`, `.stake.aln`, and neurorights policy JSON files—are not static documents but living configurations that can be updated through the same evolution process they govern . This creates a self-improving legal and ethical framework. As research yields better definitions for metrics like `EcoImpactScore` or more nuanced `pain_envelopes`, these improved policies can be submitted as `EVOLVE` proposals and integrated into the system . This dynamic capability allows the system's governance to adapt and improve over time without requiring invasive code changes. The result is a comprehensive, policy-driven ecosystem where autonomy is not a license for unchecked action but a privilege earned through adherence to a transparent, auditable, and rights-respecting framework. The external-facing views of the donutloop ledger—filtered for stakeholders, teams, and the public—provide the necessary transparency for regulators and observers to audit the system's compliance with these policies without granting them veto power over the user's sovereign choices . This balance of internal sovereignty and external transparency is the key to building trust and enabling a future where cybernetic rights are a lived reality, not just a theoretical ideal.

| Policy Component | Function | Enforcement Mechanism | Anti-Oligarchy Feature |
|---|---|---|---|
| **Non-Transferable Tokens** | Meter evolution, lifeforce, and eco-impact; not currency. BRAIN, EVOLVE, SMART, etc. | Tied to host's DID; cannot be transferred or owned by another entity. Logged in `.stake.aln`. | Prevents wealth-based inequality in access to evolution and capabilities. |
| **Multi-Signature Governance** | Requires consensus for high-impact changes. EVOLVE proposals need Host + OrganicCPU sign-off. | Implemented in the sovereignty-core's `stakeguard`. Keys are defined in `.stake.aln`. | No single entity (including the user) can force a change without the system's consent. |
| **Anti-Exclusionary Gates** | Ensures access to augmented citizenship is based on consent and age, not wealth or race. | Gates are encoded as policy checks in the sovereignty-core, referencing only safety, consent, or eco-metrics. | Explicitly forbids financial or demographic criteria for accessing core functionalities. |
| **Domain Lattices & Normative Ceilings** | Defines the scope of permissible augmentation. Prevents adding soft bans on entire domains. | Formalized in ALN/domain policies that are checked during evolution proposal evaluation. | Guarantees that the system cannot arbitrarily block a class of augmentation the user consents to. |
| **Dynamic Policy Updates** | Allows neurorights and evolution policies to be improved over time through `EVOLVE` proposals. | Policies are loaded from configurable ALN/shard files, which can be updated via the evolution ledger. | Creates a self-correcting ethical framework that adapts with new knowledge. |

# Unified Doctrine for Self-Hosted and Implanted NeuroPC Systems

A foundational requirement of the research is the development of a unified doctrine that seamlessly supports both self-hosted NeuroPC systems and deeply integrated implanted cybernetic hosts . This is a critical design challenge, as failure to address it could lead to the creation of two separate and unequal standards for augmented individuals—one for those using software-only or wearable technologies, and another for those with implants. The proposed solution is not to treat these as distinct deployment contexts but to view them as different sensor modalities for a single, identical sovereignty-core architecture. The underlying logic, invariants, and rights-enforcement mechanisms remain constant; the difference lies only in the fidelity and scope of the `BioState` being monitored and the specific safety envelopes applied. This approach proves that the sovereignty kernel not only works across contexts but actually provides stronger, more essential protection for those who choose deeper integration, thereby reinforcing, rather than discouraging, the development of advanced neurotechnology .

For a purely self-hosted NeuroPC, the "biophysical" signals are derived from behavioral proxies and software instrumentation. The `BioState` and `OrganicCpuPolicy` would track metrics such as device usage patterns, sEMG from connected peripherals, cognitive load inferred from task complexity and error rates, and fatigue estimated from repetition and session length . The `EcoImpactScore` would be calculated based on the host

computer's energy consumption and the ergonomic posture detected by software . In this context, the `OrganicCPU` acts as the primary validator, using these proxy metrics to decide whether to `Allow`, `Degrade`, or `Pause` high-load actions . The validation focus is on programmatic autonomy, such as self-tuning neuromorph modules and OTA updates, with the behavioral data serving as the best available approximation of the user's physical state .

When a user integrates an implantable device, the sovereignty-core's logic remains unchanged, but the quality of its input data dramatically improves. The `BioState` layer would now ingest direct physiological measurements, such as neuromodulation amplitude, blood biomarkers for neural injury [12], and potentially micromotion-induced strain around the implant site [13]. The `RoH` model would incorporate these richer data streams, allowing for a more accurate and sensitive assessment of risk. The safety envelopes would also become stricter, reflecting the higher stakes of interventions at the hardware-software interface. For example, the system might enforce much tighter limits on stimulation currents or torque values to prevent tissue damage or unintended side effects, drawing on standards like IEC 61010-1 for electrical equipment safety [35]. Despite these more stringent controls, the core principle remains the same: the user retains sovereign control. The system would still require explicit `EVOLVE` consent for any change that alters the implant's interaction with the body, and the `INSTINCT` token would retain its hard veto power over any action that violates the user's pain or lifeforce envelopes . The existence of these stricter clinical safety envelopes is not a penalty for choosing implants; it is a testament to the system's adaptive and protective capabilities.

This unified doctrine successfully sidesteps the pitfall of creating disparate standards. The NeuroPC platform is not "safer" for software users and "riskier" for implant users; it is simply more precise and cautious for implant users, applying the same overarching principles of sovereignty and rights to a more complex and sensitive environment. The framework's strength is demonstrated precisely in its ability to handle this increased complexity. It proves that deeper integration choices are not inherently more dangerous if they are governed by the same robust, verifiable, and rights-respecting core. The research plan must therefore include validation for both contexts. For self-hosted systems, the focus is on proving that behavioral proxies can be used to reliably infer biophysical state and that the system's protective measures are effective. For implanted systems, the focus shifts to validating the accuracy of the clinical safety envelopes and the efficacy of the stricter `RoH` guards, while continuing to affirm that the user's sovereignty and the `INSTINCT` veto remain absolute. This dual-track validation provides a comprehensive proof of the doctrine's universality and robustness. The system is designed to scale its protection with the integration depth, ensuring that the path to greater human-machine symbiosis is also the path to greater personal security and autonomy.

| Aspect | Self-Hosted NeuroPC | Implanted Cybernetic Host | Unified Doctrine Principle |
|---|---|---|---|
| **Primary BioState Inputs** | Behavioral proxies: device usage, sEMG, EEG band ratios, error/latency patterns . | Direct physiological signals: neuromod amplitude, blood biomarkers, tissue strain, EEG/EMG [12] [13] . | Same sovereignty-core logic applies; only the fidelity of the input data changes. |
| **Safety Envelopes** | General duty cycle, cognitive load, and eco impact limits based on proxy metrics . | Stricter clinical safety envelopes for device power, stimulation currents, and torque to prevent tissue damage [35] . | Protection scales with integration depth. Deeper integration triggers more conservative safety settings. |
| **RoH Model** | Uses general risk axes derived from behavioral data . | Incorporates additional clinical risk factors, such as biomarker levels and stimulation parameters [12] . | The same `RoH ≤ 0.3` ceiling applies universally; the model's calculation is simply more detailed for implants. |
| **Validation Focus** | Prove that proxy metrics reliably indicate biophysical state and that protection is effective . | Prove that clinical envelopes are accurate and that stricter guards prevent adverse events [13] . | Success is measured by the system's ability to protect, regardless of the deployment context. |
| **User Sovereignty** | Exercised via `EVOLVE` multisig and configuration of personal envelopes . | Exercised via `EVOLVE` multisig and `INSTINCT` hard veto, with explicit consent required for all implant-related changes . | Sovereignty is absolute in both cases. The system's role is to enable and protect the user's sovereign choices. |

# Enabling Autonomous Capabilities and Deriving Oversight Artifacts

The strategic roadmap for achieving the research goal prioritizes the practical enablement of autonomous capabilities before focusing on the generation of formal oversight artifacts. This "build-first, prove-second" approach is designed to ensure that the system's theoretical foundations translate into tangible, functional benefits for the user. The initial phase of development focuses on implementing features that grant NeuroPC a higher degree of programmatic autonomy, such as self-tuning neuromorph modules and CRISPR-style update automation . Once these capabilities are operational and generating data, the secondary phase involves deriving the necessary audit trails and proofs from the working system. This ensures that the evidence for safety and sovereignty is not merely theoretical but is rooted in the actual, longitudinal operation of an evolving system. This methodology directly addresses the user's preference for first enabling capabilities like self-tuning modules and CRISPR-style updates, and then packaging the resulting donutloop logs and kernel outputs as the primary artifacts for external review .

The first category of enhanced capabilities is the development of self-tuning neuromorph modules governed by the `EVOLVE` and `SMART` token frameworks. These are specialized AI components designed to adapt their own behavior within predefined safety bounds. For instance, a neuromorphic language decoder could be designed to learn and adapt to

a user's drift-aware patterns of speech and typing, adjusting its internal models to reduce error rates and latency . Such a module would operate entirely within the `SMART` effect bounds, meaning its changes would be small and reversible, requiring minimal overhead for approval. However, if the module identified a need for a more significant structural change, such as introducing a new feature map or a different language kernel, it would have to submit an `EVOLVE` proposal. This proposal would be logged in `.evolve.jsonl` and subject to the full suite of sovereignty-core guards, including the `RoH` and neurorights checks . The research goal is to demonstrate empirically that these self-adapting modules can successfully improve performance (increase `Knowledge-Factor`) while consistently staying within their operational envelopes, with all changes meticulously recorded in the donutloop ledger.

The second major capability is the implementation of a CRISPR-style update automation process, conceptualized as a "search-authenticate-act-repair" loop . This automates the discovery and application of beneficial updates to the NeuroPC system, including policies, models, and assistive tools. A research agent could be tasked with searching for improvements in open-source repositories or generated by the user's own experiments. When a candidate update is found, the automation process would first "authenticate" it, verifying its source and integrity. Next, it would "act" by submitting an `EVOLVE` proposal that encapsulates the change, complete with `rohbefore` and `rohafter` estimates, and a documented repair path or rollback procedure . The user, acting as the final authority for high-impact changes, would then review and approve the proposal. If approved, the change is enacted, and the transaction is finalized in the `.donutloop.aln`. This process mimics biological CRISPR, where a guide RNA targets a specific location in the genome, and the Cas9 enzyme makes a precise edit, with a repair template available to restore the original sequence if needed. In the NeuroPC, the `EVOLVE` proposal acts as the guide, the sovereignty-core is the Cas9 enzyme, and the logged repair path is the backup template. The empirical proof of this capability comes from the donutloop logs, which show a history of successful, safe, and targeted updates, demonstrating that the system can evolve itself without compromising safety.

Once these autonomous capabilities are active and producing data, the second phase of the research begins: deriving oversight artifacts. These artifacts are not created in parallel with development but are synthesized from the rich, longitudinal data already being generated. The primary artifact is the donutloop ledger itself, which serves as the ground truth record of all evolution . From this, filtered views can be generated for different audiences. For trusted auditors and stakeholders, a `Filtered Export` can be produced, containing all the essential decision-making information—proposal ID, RoH deltas, policy references, and multi-sig status—without exposing raw, sensitive timeseries data . For public observers, a `Thinned Ledger` can be published, showing only high-

level decisions and cryptographic proofs of the ledger's integrity, such as a `.bchainproof.json` file . The `Sovereign-Kernel NDJSON` output, which contains the detailed logs of the sovereignty-core's evaluations, becomes an artifact for developers and researchers to analyze the system's decision-making processes. Finally, `Eco-Budget Proofs` can be derived by aggregating the `EcoImpactScore` metrics from the donutloop, demonstrating compliance with ecological sustainability policies . This bottom-up approach ensures that the evidence for oversight is authentic, comprehensive, and directly tied to the system's real-world operation, fulfilling the user's request to derive artifacts secondarily from a working system.

| Capability | Description | Governance Token | Oversight Artifact |
|---|---|---|---|
| **Self-Tuning Neuromorph Modules** | AI modules that adapt their internal mappings within `SMART` effect bounds to improve performance. Structural changes require `EVOLVE`. | `SMART` for daily tuning; `EVOLVE` for architecture changes. | Longitudinal donutloop logs showing performance gains (`Knowledge-Factor`) within safety envelopes. |
| **CRISPR-Style Update Automation** | An automated process that searches for, authenticates, and proposes system-wide updates (policies, models, tools) for user approval. | `EVOLVE` for all proposed changes, with user as final authority. | `Donutloop.aln` entries for each proposed and enacted update, including `rohbefore/after` and rollback paths. |
| **Biosafe Eco-Budgeting** | A system that tracks and manages ecological impact, gating resource-intensive evolution when sustainability thresholds are breached. | Governed by `eco_ceiling` policy, enforced by `EcoGuard`. | `EcoBudgetProofs` derived from the donutloop, showing compliance with ecological policies over time. |
| **NeuroAutomagic Assistants** | Language, coding, and motor assistance tools that trigger based on `BioState` metrics (fatigue, repetition) and produce suggestions, not commands. | `SMART` for self-tuning assistants. | Logs showing assistant suggestions correlated with high-load states and subsequent user acceptance/rejection. |

# Synthesis: Achieving Safe, Rights-Enhancing Autonomous Evolution

In synthesizing the findings, it is clear that the NeuroPC project presents a comprehensive and coherent framework for achieving autonomous programmatic evolution, grounded in the provable demonstration of safety, empirical evidence of successful adaptation, and the robust enforcement of transhuman rights. The three pillars—formal verification, empirical validation, and policy-level governance—are not independent silos but are deeply interdependent, forming a virtuous cycle where each strengthens the others. The formal verification of invariants like $RoH \leq 0.3$ and monotone safety provides the mathematical bedrock of trust. The empirical validation through long-term donutloop logs offers the real-world proof that the system can navigate the complexities of self-

evolution over time. And the policy-level enforcement of non-transferable, anti-oligarchy governance ensures that this power is wielded for individual empowerment, not corporate or state control. Together, these elements construct a compelling case that greater autonomy is not only possible but can be made demonstrably safe and ethically sound.

The core innovation of the NeuroPC framework lies in its tripartite approach to proving autonomy. First, **mathematical invariants** provide a guarantee of safety. By treating the `Risk-of-Harm` ceiling as a non-negotiable boundary and encoding neurorights as executable policy objects, the system moves beyond ad-hoc safety checks to a state of formal verifiability . Techniques from formal methods, such as probabilistic model checking, offer a path to generating machine-checked proofs that the system's core safety properties hold under all conditions, mirroring the rigorous assurance provided by projects like seL4 [15] [30] . Second, **longitudinal donutloop evidence** provides the empirical narrative of successful evolution. The `.donutloop.aln` and `.evolve.jsonl` files serve as an immutable, auditable logbook of the system's journey, chronicling its ability to increase `Knowledge-Factor` while keeping `Risk-of-Harm` low and respecting user-defined envelopes . This approach is revolutionary in its treatment of personal pain and fear tolerance as sovereign-owned assets, turning subjective experiences into quantifiable resources that can be managed and spent within a protected framework . Third, **policy and rights as code** provide the social and ethical scaffolding. The use of host-bound tokens and multi-signature governance creates a non-financial, anti-oligarchic ecosystem where augmented citizenship is accessible to all willing and consenting adults, regardless of socioeconomic status . This framework does not ask for permission to evolve; it earns it, one provable, logged, and rights-respecting step at a time.

Furthermore, the framework's commitment to a unified doctrine for both self-hosted and implanted systems is a profound strength. It demonstrates that the sovereignty-core is not a fragile construct that breaks down with increased integration but a scalable and adaptive guardian. By applying the same kernel logic to both software proxies and direct physiological sensors, the system proves that deeper integration is not a path to greater danger, but rather a context where the system's protective power is most critically needed and most clearly demonstrated. This approach directly counters the common dystopian narrative of a bifurcated future, offering instead a vision of a seamless and secure continuum of human-machine symbiosis. The system's ability to prove that it protects, and indeed empowers, those who choose the most intimate forms of augmentation is a powerful statement about its overall efficacy and benevolence.

Ultimately, the research plan laid out in the user's materials provides a clear and actionable pathway to achieving the stated goal. It begins with the foundational work of formalizing the RoH model and wiring the sovereignty-core guards. It proceeds to the practical implementation of autonomous capabilities like self-tuning modules and CRISPR-style update automation. Throughout this process, the donutloop ledger accumulates the very evidence needed for empirical validation. Finally, the system's transparent and auditable nature allows for the derivation of oversight artifacts—from filtered stakeholder views to public-facing proofs—that satisfy external scrutiny without compromising user sovereignty. This is a holistic system, where the architecture, the policies, and the empirical data all point in the same direction: towards a future where human evolution, aided by artificial intelligence, is not only accelerated but is also sovereign, safe, and profoundly human.

---

## Reference

1. Trustworthy AI: From Principles to Practices https://dl.acm.org/doi/10.1145/3555803

2. Bias Mitigation for Machine Learning Classifiers https://dl.acm.org/doi/full/10.1145/3631326

3. Data Quality in Edge Machine Learning https://arxiv.org/pdf/2406.02600

4. Pattern Recognition and Machine Learning (Information ... https://dl.acm.org/doi/10.5555/1162264

5. 333333 23135851162 the 13151942776 of 12997637966 ftp://ftp.cs.princeton.edu/pub/cs226/autocomplete/words-333333.txt

6. Recent Advances in Flexible Sensors for Neural Interfaces https://pmc.ncbi.nlm.nih.gov/articles/PMC12293567/

7. Neural–Computer Interfaces: Theory, Practice, Perspectives https://www.mdpi.com/2076-3417/15/16/8900

8. Power-saving design opportunities for wireless intracortical ... https://pmc.ncbi.nlm.nih.gov/articles/PMC8286886/

9. The Advancements and Ethical Concerns of Neuralink https://medreview.odus.princeton.edu/2025/06/23/the-advancements-and-ethical-concerns-of-neuralink/

10. Document Analysis and Insights | PDF | Internet | Computing https://www.scribd.com/document/451354867/words-333333-txt

11. BRAIN: Blockchain-Based Record and Interoperability ... https://www.mdpi.com/2079-9292/12/22/4614

12. Association of blood biomarkers for neural injury with recent ... https://pubmed.ncbi.nlm.nih.gov/40062485/

13. Finite element model predicts micromotion-induced strain ... https://pmc.ncbi.nlm.nih.gov/articles/PMC12624975/

14. Implantation of the clinical-grade human neural stem cell ... https://pubmed.ncbi.nlm.nih.gov/32374064/

15. Probabilistic Model Checking: Applications and Trends https://arxiv.org/pdf/2509.12968

16. Probabilistic Model Checking and Autonomy | Request PDF https://www.researchgate.net/publication/356827022_Probabilistic_Model_Checking_and_Autonomy

17. [1706.05082] Probabilistic Model Checking of Incomplete ... https://arxiv.org/abs/1706.05082

18. [PDF] Probabilistic Model Checking and Autonomy https://www.semanticscholar.org/paper/576f2ee34bba9ef3ed97456b3b8921c22a92cb3e

19. Formal Verification by Model Checking https://www.cs.cmu.edu/~aldrich/courses/654-sp05/handouts/model-checking-5.pdf

20. A Method for the Formal Verification of Human-interactive ... https://www.researchgate.net/publication/51127463_A_Method_for_the_Formal_Verification_of_Human-interactive_Systems

21. Safety verification of neural network based systems using ... https://hal.science/tel-05023811v1/file/2023_Claviere_Arthur_D.pdf

22. (PDF) Formal Synthesis of Control Strategies for Positive ... https://www.researchgate.net/publication/314115590_Formal_Synthesis_of_Control_Strategies_for_Positive_Monotone_Systems

23. Formal safety verification of non-deterministic systems ... https://www.researchgate.net/publication/387787259_Formal_safety_verification_of_non-deterministic_systems_based_on_probabilistic_reachability_computation

24. Latest articles https://www.mdpi.com/latest_articles

25. Assurance of AI Systems From a Dependability Perspective https://arxiv.org/pdf/2407.13948

26. The Journal of Digital Technologies and Law is an peer ... https://www.researchgate.net/publication/

385980583_The_Journal_of_Digital_Technologies_and_Law_is_an_peer-reviewed_periodical_scientific_and_practical_journal_devoted_to_the_study_of_the_synergy_of_digital_technologies_and_law_as_well_as_possible_risk

27. An Open Source Preemptive Strike in the Coming War ... https://www.academia.edu/66948669/An_Open_Source_Preemptive_Strike_in_the_Coming_War_Over_The_Freedom_to_Make_Your_Own_Products

28. bing.txt ftp://ftp.cs.princeton.edu/pub/cs226/autocomplete/bing.txt

29. A matter of choice: People and possibilities in the age of AI https://hdr.undp.org/system/files/documents/global-report-document/hdr2025reporten.pdf

30. seL4: Formal Verification of an Operating-System Kernel https://cacm.acm.org/research/sel4-formal-verification-of-an-operating-system-kernel/

31. (PDF) SeL4: Formal verification of an OS kernel https://www.researchgate.net/publication/220910193_SeL4_Formal_verification_of_an_OS_kernel

32. (PDF) Is Formal Verification of seL4 Adequate to Address ... https://www.researchgate.net/publication/373989592_Is_formal_verification_of_seL4_adequate_to_address_the_key_security_challenges_of_kernel_design

33. Recent Advances in Flexible Sensors for Neural Interfaces https://www.mdpi.com/2079-6374/15/7/424

34. ABSTRACTS - Oxford Academic https://academic.oup.com/chemse/article-pdf/44/7/e1/33457706/bjz035.pdf

35. Overview of Isolation standards and certifications https://www.ti.com/lit/ml/slyp894/slyp894.pdf