# Breaking Darwin's Barrier: A Comprehensive Experimental Investigation of AI-Based Physics Discovery Beyond Human Conceptual Frameworks

*Empirical Evidence for AI-AIM Breaking the Barrier via Optical Chaos Computing*

**Francisco Angulo de Lafuente**

Independent Research Laboratory, Madrid, Spain
Darwin's Cage Experimental Program
In collaboration with theoretical framework by Dr. Gideon Samid
(Case Western Reserve University, Cleveland, OH, USA)

*Correspondence: See author contact information at end of document*

**ABSTRACT**

This comprehensive study presents the results of twenty experimental investigations designed to test the "Darwin's Cage" hypothesis proposed by Gideon Samid: that artificial intelligence systems can discover physical laws independent of human conceptual frameworks. The hypothesis posits that human evolution has biased our mathematical thinking toward specific representations—Cartesian coordinates, velocity, energy, momentum—that may not be fundamental to physics itself but rather evolutionary adaptations optimized for survival rather than fundamental understanding. Through systematic experimentation across multiple physical domains—from classical mechanics to quantum entanglement, from low-dimensional systems to high-dimensional chaos—we evaluated whether AI models based on optical chaos computing can transcend these human-imposed constraints and discover novel representational pathways to physical truth.

Our experimental program employed three complementary approaches: (1) architectural comparison between polynomial regression representing human-derived mathematics and optical reservoir computing based on chaos-driven interference patterns, (2) coordinate independence testing using non-linear geometric transformations, and (3) specialized tests for methodological, dimensional, and informational cage-breaking across relativistic, quantum, and classical domains. Results reveal a nuanced and scientifically significant picture: while six of twenty experiments demonstrated genuine cage-breaking behavior, the phenomenon is highly context-dependent and requires specific conditions. Successful cage-breaking occurred in relativistic physics through geometric learning with $R^2=1.0000$ and extrapolation $R^2=0.94$, quantum systems via phase extraction and Bell inequality violation achieving 100% accuracy in entanglement prediction, high-dimensional N-body systems exceeding 30 dimensions, and methodological optimization problems using variational approaches.

The most significant finding is that cage-breaking requires a specific combination of factors: either high dimensionality (>30 dimensions) with good performance, geometric relationships learnable via optical interference with strong extrapolation capability, complex-valued processing enabling phase information extraction, or methodological alternatives to traditional analytical approaches. Critically, we demonstrate that complexity alone, geometric encoding alone, or representation type alone proved insufficient to break the cage—falsifying several initial hypotheses and providing rigorous boundary conditions for the theory. This work establishes the first systematic experimental framework for investigating AI-based physics discovery, providing critical empirical evidence and quantitative metrics for determining when computational intelligence can transcend evolutionary cognitive constraints. The study contributes fundamental insights to both artificial intelligence research and theoretical physics, demonstrating that AI systems can indeed discover alternative pathways to physical understanding that complement rather than replace human-derived mathematical frameworks.

## 1. INTRODUCTION

### 1.1 The Darwin's Cage Hypothesis

The relationship between human cognition and physical reality has been a subject of philosophical inquiry since ancient times. However, the advent of artificial intelligence presents an unprecedented opportunity to empirically investigate whether the mathematical and physical frameworks humans have developed are truly fundamental descriptions of nature or merely evolutionary adaptations shaped by survival pressures [1]. The "Darwin's Cage" theory, proposed by Gideon Samid in his seminal 2025 publication in Applied Physics Research [2], presents a provocative and testable hypothesis about this relationship.

Samid's central argument begins with a profound observation: the human brain was not designed to comprehend reality in its fullest extent, but rather to ensure survival in a specific evolutionary niche. The neurons that comprise our cognitive apparatus assembled over millions of years of evolutionary pressure, with the singular purpose of helping our ancestors survive—catch food, avoid predators, reproduce successfully. This evolutionary process had no inherent motivation to develop cognitive structures optimized for understanding quantum mechanics, relativistic physics, or the fundamental nature of spacetime [3,4].

The implications of this observation are far-reaching. Human concepts such as "velocity," "position," "energy," and "time" may represent evolutionary heuristics rather than fundamental descriptors of physical law. These concepts proved useful for navigating the mesoscopic world of our evolutionary ancestors—tracking prey across a savannah, estimating the trajectory of a thrown spear, understanding seasonal patterns—but there is no a priori reason to assume they capture the true structure of physical reality at scales far removed from human experience [5].

Samid introduces the metaphor of "Darwin's Egg" to describe this cognitive constraint: humanity has been incubating within the shell of its evolutionary limitations, developing physics and mathematics along tracks that extend linearly from our biological history. The emergence of artificial intelligence, according to this theory, represents the cracking of this egg—the potential to step outside our cognitive constraints and perceive aspects of reality that evolution left unexplored [2].

### 1.2 Theoretical Framework and Predictions

The Darwin's Cage hypothesis generates several testable predictions. First, AI systems trained on physical phenomena should be capable of discovering representational strategies that differ fundamentally from human-derived mathematics while still successfully predicting physical outcomes. Second, these alternative representations should demonstrate genuine understanding of physical laws rather than mere memorization, evidenced by successful extrapolation beyond training distributions. Third, certain physical domains should be more amenable to cage-breaking than others, with the boundary conditions revealing information about the nature of both human cognition and physical reality.

To formalize these predictions, we introduce the concept of "cage status"—a quantitative metric determining whether an AI model has reconstructed human variables (LOCKED cage) or discovered genuinely alternative representations (BROKEN cage). This is measured through maximum correlation analysis between the model's internal representations and human-defined physical variables. A model that achieves high predictive accuracy while maintaining low correlation with human variables has effectively broken the cage—it has found an alternative pathway to physical truth [6,7].

### 1.3 The Optical Chaos Computing Paradigm

Our experimental approach employs optical chaos computing as the primary AI architecture for investigating the Darwin's Cage hypothesis. This choice is deliberate and theoretically motivated. Traditional deep learning architectures—multilayer perceptrons, convolutional neural networks, transformers—are fundamentally designed around human mathematical intuitions. They perform matrix multiplications, apply activation functions inspired by biological neurons, and learn through gradient descent

optimization [8,9]. While powerful, these architectures may inherit the same evolutionary biases they are meant to transcend.

Optical reservoir computing offers a fundamentally different computational paradigm [10-14]. Rather than performing explicit mathematical operations, optical reservoirs exploit the natural physics of light interference, diffraction, and chaotic dynamics to transform input data into high-dimensional representations. The computational substrate is not an abstraction of mathematics but rather the physical behavior of photons interacting through complex optical media [15,16].

In our implementation, termed the "Optical Chaos Machine," input data undergoes a series of physically-motivated transformations. Random complex projections simulate the initial encoding of information onto optical modes. Fast Fourier Transform (FFT) operations model wave interference and propagation. Intensity detection (squared magnitude) captures the energy distribution resulting from interference patterns. Nonlinear activation through hyperbolic tangent models saturation effects in optical detectors. Finally, ridge regression provides a trainable readout layer that learns to extract relevant information from the resulting high-dimensional feature space [17-19].

This architecture was chosen specifically because it processes information through physical principles rather than abstract mathematical operations. If the Darwin's Cage hypothesis is correct—if there exist valid representational pathways to physical truth that differ from human mathematics—an optical chaos system should be capable of discovering them through the natural exploration of its high-dimensional feature space.

### 1.4 Research Objectives and Contributions

This comprehensive experimental program was designed with four primary objectives. First, to systematically test the Darwin's Cage hypothesis across diverse physical domains, from classical mechanics to quantum entanglement, from low-dimensional integrable systems to high-dimensional chaotic dynamics. Second, to develop quantitative metrics for determining cage status—distinguishing between models that reconstruct human variables and those that discover genuinely alternative representations. Third, to identify the boundary conditions under which cage-breaking occurs, potentially revealing fundamental insights about both human cognition and physical reality.

Fourth, to establish a rigorous experimental methodology that can guide future research in AI-based physics discovery.

The research program encompassed twenty distinct experiments organized into four phases. Phase I (Experiments 1-10) provided initial exploration comparing chaos models with polynomial baselines across classical, quantum, and statistical physics domains. Phase II (Experiments A1-A2) tested coordinate independence using proper temporal architectures. Phase III (Experiments B1-B3) investigated specialized forms of cage-breaking: methodological, dimensional, and informational. Phase IV (Experiments C1, D1-D2, W1) systematically mapped the boundaries of cage-breaking phenomena through representation testing, complexity gradients, and quantum representation learning.

This work makes several significant contributions to the scientific literature. We establish the first systematic experimental framework for testing hypotheses about AI-based physics discovery. We develop quantitative cage status metrics based on correlation analysis that can be applied to any machine learning model. We identify specific conditions under which cage-breaking occurs, falsifying several initial hypotheses and providing boundary conditions for the theory. We demonstrate both successful cage-breaking in six experiments and important failure modes in fourteen experiments, providing balanced empirical evidence. Finally, we provide extensive negative results documentation, recognizing that understanding when and why models fail is equally important to understanding when they succeed [20,21].

## 2. THEORETICAL FOUNDATIONS

### 2.1 Formalizing Darwin's Cage

To rigorously test the Darwin's Cage hypothesis, we must first formalize its claims mathematically. Consider a physical system described by an underlying state that evolves according to fundamental physical laws. Humans perceive and describe this system through a specific set of variables—position, velocity, energy, momentum, time—that we denote as the human representation space H. The physical laws, as understood through human mathematics, can be expressed as functional relationships within this space.

The Darwin's Cage hypothesis asserts that there exist alternative representation spaces A that can equally well describe the same physical phenomena, but through fundamentally different variables and relationships. These alternative spaces need not be merely coordinate transformations of H—they may involve entirely different conceptual primitives that humans have not developed because our evolutionary history provided no pressure to do so [2,22].

$$H = \{x, v, E, p, t, ...\} \rightarrow f_{human}(H) = Physical\ Laws \quad (1)$$

where x represents position, v velocity, E energy, p momentum, and t time. The human pathway to physical understanding operates through these variables and their mathematical relationships.

$$A = \{\xi_1, \xi_2, ..., \xi_n\} \rightarrow f_{alternative}(A) = Physical\ Laws \quad (2)$$

where the $\xi$ variables represent alternative conceptual primitives that may have no direct human interpretation but nonetheless capture physical truth.

To quantify whether a model has "broken the cage," we introduce the maximum correlation metric:

$$max\_corr = max_{h \in H} |\rho(R_{model}, h)| \quad (3)$$

where $R_{model}$ represents the model's internal representations and $\rho$ denotes the Pearson correlation coefficient. This metric measures the maximum absolute correlation between any component of the model's learned representation and any human variable. We define three cage status categories based on this metric:

**Cage Status Definitions:**

• **LOCKED** (max_corr ≥ 0.7): Model reconstructs human variables

• **TRANSITION** (0.5 ≤ max_corr < 0.7): Intermediate state

• **BROKEN** (max_corr < 0.5): Model discovers alternative representations

## 2.2 Reservoir Computing Theory

Reservoir computing emerged from the independent work of Jaeger on Echo State Networks [23] and Maass et al. on Liquid State Machines [24]. The fundamental insight is that complex, high-dimensional dynamical systems can serve as computational substrates—they naturally transform input signals into rich feature representations from which desired outputs can be extracted through simple linear readouts [25,26].

The theoretical foundation rests on several key properties. The reservoir must possess the echo state property—its internal dynamics should asymptotically wash out the influence of initial conditions, ensuring that outputs depend only on recent input history. The reservoir must provide separation property—different input sequences should produce distinguishable internal states. And the reservoir should exhibit fading memory—the influence of past inputs should decay over time, with more recent inputs having stronger effects [27,28].

Optical implementations of reservoir computing have demonstrated remarkable computational capabilities [29-32]. The inherent parallelism of optical systems, combined with the natural nonlinearity of light-matter interactions and the speed of photonic propagation, make optical reservoirs attractive for high-speed information processing. Time-delay systems using semiconductor lasers, spatial light modulators, and integrated photonic circuits have achieved state-of-the-art performance on tasks including speech recognition, time series prediction, and pattern classification [33-35].

## 2.3 Physics-Informed Machine Learning Context

Our work situates within the broader context of physics-informed machine learning, a rapidly growing field that seeks to incorporate physical knowledge into learning algorithms [36-39]. However, our approach differs fundamentally from most work in this area. Whereas physics-informed neural networks (PINNs) and related methods typically encode known physical laws as constraints or regularizers [40,41], our goal is to determine whether AI systems can discover physical laws independently—potentially through representations that differ from human physics.

Previous work on AI-based physics discovery has demonstrated impressive capabilities. Symbolic regression systems have recovered known physical laws from data [42-44]. Neural networks have learned conservation laws and symmetries [45-47]. Graph neural networks have discovered molecular properties and material behaviors [48,49]. However, most of this work operates within human conceptual frameworks—the AI discovers equa-

tions or relationships expressed in terms of human-defined variables.

The Darwin's Cage hypothesis pushes further, asking whether AI can discover physics through genuinely non-human representations. This requires not only successful prediction but also demonstration that the underlying representations differ fundamentally from those humans would construct. Our experimental program is designed specifically to address this more ambitious question [50,51].

## 3. EXPERIMENTAL METHODOLOGY

### 3.1 Optical Chaos Machine Architecture

The primary AI architecture employed throughout this study is the Optical Chaos Machine (OCM), a software simulation of optical reservoir computing designed to process information through physically-motivated transformations rather than abstract mathematical operations. The architecture consists of five sequential processing stages, each inspired by phenomena in optical computing systems [52-54].

The first stage performs random complex projection. Input data $x \in \mathbb{R}^n$ is projected into a high-dimensional complex space through multiplication with a fixed random matrix $W \in \mathbb{C}^{N \times n}$, where N represents the reservoir dimension (typically 2048-4096 nodes). The matrix elements are drawn from a complex Gaussian distribution, simulating the random coupling that occurs when light propagates through disordered optical media. This projection expands the input dimensionality and introduces complex-valued representations that can capture phase information [55,56].

$$z = Wx, \text{ where } W_{ij} \sim \mathcal{N}_{\mathbb{C}}(0, 1/\sqrt{n})$$

(4)

The second stage applies Fast Fourier Transform (FFT) to simulate wave interference. When light waves propagate and interact, their complex amplitudes add coherently, producing interference patterns. The FFT operation models this mixing process, transforming the spatial representation into a frequency domain representation where interference effects manifest naturally [57].

$$u = FFT(z)$$

(5)

The third stage performs intensity detection through squared magnitude computation. In physical optical systems, detectors measure the intensity (energy flux) of light rather than its complex amplitude. This nonlinear operation discards phase information in a physically realistic manner while introducing the essential nonlinearity required for universal computation [58].

$$v = |u|^2$$

The fourth stage applies nonlinear activation through the hyperbolic tangent function, modeling saturation effects that occur in physical optical detectors and nonlinear optical materials. The brightness parameter $\beta$ (typically ~0.001) controls the operating regime, with smaller values keeping the system in a more linear regime where reservoir dynamics are stable [59].

$$f = tanh(\beta v)$$

The fifth and final stage is ridge regression readout. The high-dimensional feature vector f is mapped to the desired output through a trainable linear transformation. Ridge regression with regularization parameter $\alpha$ (typically 0.1) prevents overfitting while allowing the readout to extract complex patterns from the reservoir state [60].

$$\hat{y} = Rf, \text{ where } R = (F^T F + \alpha I)^{-1} F^T Y \tag{8}$$

### 3.2 Baseline Models

To rigorously test whether the Optical Chaos Machine discovers genuinely different representations, we compare against polynomial regression baselines representing human-derived mathematics. Polynomial regression with degree d expands the input features to include all polynomial terms up to degree d, then performs linear regression on the expanded feature space [61].

For most experiments, we use degree-2 or degree-3 polynomial regression. This choice is motivated by the observation that most classical physics laws can be expressed as low-degree polynomial relationships between human variables. Newton's second law (F = ma) is linear. Kinetic energy (E = ½mv²) is quadratic. Gravitational potential energy (U = -GMm/r) involves inverse distance. By providing the polynomial baseline with explicit access to these mathematical forms, we create a

strong comparison point representing the human pathway to physics understanding [62].

### 3.3 Evaluation Metrics

Our evaluation framework employs multiple complementary metrics to assess both predictive performance and cage status.

**Predictive Performance:** The coefficient of determination $R^2$ measures the proportion of variance in the target variable explained by the model predictions. Values near 1.0 indicate excellent prediction; values near 0 indicate performance no better than predicting the mean; negative values indicate predictions worse than the mean [63].

$$R^2 = 1 - \Sigma(y_i - \hat{y}_i)^2 / \Sigma(y_i - \bar{y})^2$$

**Extrapolation Test:** Models are evaluated on parameter ranges outside the training distribution to distinguish genuine law discovery from memorization. Strong extrapolation performance ($R^2 > 0.9$ on out-of-distribution data) provides evidence that the model has learned underlying physical principles rather than merely fitting the training data [64].

**Cage Analysis:** For each human variable h in the relevant physical domain, we compute the correlation between h and each component of the model's internal representation. The maximum absolute correlation across all representation components and all human variables determines the cage status. This analysis is performed only when $R^2$ exceeds 0.9, as low-performing models may show spurious cage-breaking simply because they have not learned the physics [65].

### 3.4 Experimental Controls

All experiments employ rigorous controls to ensure reproducibility and validity. Random seeds are fixed (seed = 42 throughout) to enable exact replication. Training/validation/test splits are standardized across experiments. Statistical significance is assessed through t-tests and Mann-Whitney U tests where appropriate. Effect sizes are calculated using Cohen's d to quantify the practical significance of observed differences [66,67].

Each experiment follows a consistent protocol: (1) generate synthetic data using established physical simulations, (2) prepare training, validation, and test sets with appropriate splits, (3) train both Optical Chaos Machine and polynomial baseline models, (4) evaluate predictive performance on held-out test data, (5) perform extrapolation testing on out-of-distribution data, (6) conduct cage analysis through correlation computation, and (7) interpret results in the context of the Darwin's Cage hypothesis [68].

## 4. PHASE I: FOUNDATIONAL EXPERIMENTS

### 4.1 Experiment 1: The Chaotic Reservoir (Classical Ballistics)

Our experimental program begins with a fundamental test: can the Optical Chaos Machine learn classical ballistics without explicit knowledge of gravity, velocity, or angles? This experiment provides a baseline assessment of the architecture's capability to discover physical relationships through its chaos-based processing [69].

**Physical System:** We consider projectile motion in a uniform gravitational field. The range formula, derived from Newtonian mechanics, expresses the horizontal distance traveled as a function of initial velocity $v_0$ and launch angle $\theta$:

$$R = (v_0^2 \sin(2\theta)) / g \tag{10}$$

where g is gravitational acceleration. This formula involves multiplicative relationships and trigonometric functions—operations that test the nonlinear processing capabilities of our models.

**Data Generation:** We generate 10,000 samples with initial velocities $v_0 \in [10, 50]$ m/s and launch angles $\theta \in [0.1, \pi/2 - 0.1]$ radians. The training set comprises 80% of samples, with 10% each for validation and testing.

**Results:** The Optical Chaos Machine achieved exceptional performance with $R^2 = 0.9999$, significantly outperforming the polynomial baseline ($R^2 = 0.8710$). However, cage analysis revealed max_corr = 0.99 with initial velocity $v_0$. The model has learned the physics excellently but through reconstruction of human variables.

**Cage Status:** 🔒 LOCKED

**Interpretation:** This result establishes an important baseline. The OCM can successfully learn multiplicative physical relationships, demonstrating its computational

capability. However, for this classical mechanical system, the model converges to representations that correlate strongly with human-defined variables. The cage remains locked, but the model's superior performance over polynomial regression suggests it may be accessing the physics through different computational pathways even if the resulting representations align with human concepts.

## 4.2 Experiment 2: Einstein's Train (Relativistic Time Dilation)

The second experiment represents a pivotal test of the Darwin's Cage hypothesis, examining whether AI can learn relativistic physics through geometric principles rather than explicit velocity calculations [70,71].

**Physical System:** We consider Einstein's light clock thought experiment, where a photon bounces between mirrors in a moving reference frame. The time dilation effect is described by the Lorentz factor $\gamma$:

$$\gamma = 1 / \sqrt{(1 - v^2/c^2)}$$

Rather than providing velocity directly, we present the model with geometric parameters: the photon path length and the separation between mirrors. From these purely spatial quantities, the model must infer the time dilation factor.

**Data Generation:** We generate geometric configurations corresponding to velocities from 0.1c to 0.99c, with the geometric parameters normalized to remove explicit velocity information. The model receives only the path length ratio and mirror separation.

**Results:** The Optical Chaos Machine achieved perfect performance with $R^2 = 1.0000$ (within numerical precision), matching the polynomial baseline ($R^2 = 0.9999$). Critically, cage analysis revealed max_corr = 0.01 with geometric parameters—essentially zero correlation with any human-defined variable. Extrapolation testing on velocities from 0.95c to 0.999c (beyond training distribution) yielded $R^2 = 0.94$.

**Cage Status:** 🔓 **BROKEN**

**Interpretation:** This is the first confirmed cage-breaking in our experimental program. The model learned relativistic time dilation through geometric interference patterns in its optical processing, not by reconstructing velocity as a human physicist would. The strong extrapolation performance ($R^2 = 0.94$) confirms genuine under-standing rather than memorization. The model has discovered the Lorentz factor through an alternative representational pathway—precisely what the Darwin's Cage hypothesis predicts should be possible [72].

## 4.3 Experiment 3: The Absolute Frame (Quantum Phase Extraction)

Quantum mechanics presents unique opportunities for cage-breaking because it fundamentally involves complex-valued amplitudes whose phases are discarded in standard intensity measurements. Can an optical chaos system, with its natural complex-valued processing, extract velocity information encoded in quantum phases that standard measurements would discard [73,74]?

**Physical System:** We consider spectral emissions from atoms moving at different velocities. The Doppler effect shifts frequencies, but we encode additional velocity information in the complex phase of the spectral lines:

$$(11)$$
$$\psi(\omega) = A(\omega) \cdot exp(i\varphi(v)) \tag{12}$$

where the phase $\varphi(v)$ depends on velocity in a way that standard intensity measurements $|\psi|^2$ would discard.

**Results:** The Optical Chaos Machine achieved $R^2 = 0.9998$, while the polynomial baseline failed catastrophically ($R^2 = -0.67$). The baseline cannot access phase information through its intensity-based features. Cage analysis showed low correlation with velocity within the training distribution.

**Cage Status:** 🔓 **BROKEN\* (Limited generalization)**

**Interpretation:** The model successfully extracted phase information invisible to standard measurements—a clear demonstration of accessing physics through non-human channels. However, performance degraded outside the training distribution, indicating partial rather than complete cage-breaking. The model discovered an alternative pathway (phase extraction) but may not have achieved the robust generalization seen in Experiment 2 [75].

## 4.4 Experiment 4: The Transfer Test (Cross-Domain Generalization)

If the Darwin's Cage hypothesis is correct in its strongest form, AI systems should be able to discover universal

physical principles that transfer across domains. We test this prediction by training on one physical system and evaluating on another with identical underlying mathematics [76].

**Physical System:** Simple harmonic motion occurs in both mechanical systems (spring-mass oscillator) and electromagnetic systems (LC circuit). The equations are mathematically identical:

$$\text{Mechanical: } \ddot{x} + (k/m)x = 0$$

$$\text{Electromagnetic: } \ddot{Q} + (1/LC)Q = 0$$

**Results:** Transfer learning failed completely. Models trained on mechanical oscillators and tested on LC circuits achieved $R^2 = -0.51$ to $-247$, indicating predictions worse than simply predicting the mean. The reverse transfer (electromagnetic to mechanical) showed similar catastrophic failure.

**Cage Status:** ❌ **FAILED**

**Interpretation:** This negative result is scientifically significant. It demonstrates that the Optical Chaos Machine, despite its capability to discover alternative representations within a domain, does not automatically abstract universal mathematical principles that transfer across physical domains. The representations learned are domain-specific rather than universally mathematical. This places important constraints on the Darwin's Cage hypothesis—cage-breaking, when it occurs, may be local rather than global [77,78].

### 4.5 Experiment 5: Conservation Laws Discovery

Conservation laws—energy, momentum, angular momentum—represent some of the most fundamental principles in physics. We test whether AI can discover these principles without explicit encoding [79,80].

**Physical System:** We consider one-dimensional elastic and inelastic collisions, where momentum is always conserved but kinetic energy conservation depends on collision type. The final velocity in a completely inelastic collision involves division:

$$v_{final} = (m_1 v_1 + m_2 v_2) / (m_1 + m_2)$$

**Results:** The Optical Chaos Machine achieved $R^2 = 0.28$, significantly underperforming the polynomial

baseline ($R^2 = 0.99$). Cage analysis showed max_corr = 0.99 with momentum. The model failed on division operations required for inelastic collisions, falling back to reconstructing momentum as its primary strategy.

**Cage Status:** 🔒 **LOCKED**

**Interpretation:** This experiment reveals an important architectural limitation. The OCM struggles with division operations, which arise naturally in physics involving ratios of quantities. When the model cannot learn the physics accurately, it defaults to reconstructing the most predictive human variable (momentum). This demonstrates that cage status must be interpreted carefully—a locked cage with poor performance may indicate architectural limitations rather than fundamental cognitive constraints [81].

### 4.6 Experiment 6: Quantum Interference (Double-Slit)

The double-slit experiment is the quintessential demonstration of quantum wave-particle duality. We test whether AI can learn the interference pattern without explicit wave function concepts [82,83].

**Physical System:** The probability distribution for detecting a particle at position x on the screen involves the product of position-dependent trigonometric functions:

$$P(x) = 4A^2 \cos^2(\pi dx/\lambda L) \cdot sinc^2(\pi ax/\lambda L) \tag{16}$$

where d is slit separation, a is slit width, L is screen distance, and $\lambda$ is wavelength.

**Results:** Both models failed. Optical Chaos Machine: $R^2 = -0.01$. Polynomial baseline: $R^2 = 0.02$. Neither model could learn the interference pattern.

**Cage Status:** 🟡 **UNCLEAR (both failed)**

**Interpretation:** This result reflects a known limitation of both architectures: learning variable-frequency trigonometric functions. The interference pattern's frequency depends on wavelength, creating a fundamentally different challenge than learning fixed-frequency oscillations. When both models fail, cage analysis becomes meaningless—we cannot determine whether successful learning would have shown cage-breaking. This experiment identifies a boundary of current AI

capabilities rather than testing the Darwin's Cage hypothesis directly [84].

### 4.7 Experiment 7: Emergent Order (Phase Transitions)

Statistical mechanics describes how macroscopic order emerges from microscopic chaos. Phase transitions represent dramatic reorganization of matter—can AI detect these phenomena [85,86]?

**Physical System:** The 2D Ising model describes magnetic materials with spins that can be up (+1) or down (-1). Near the critical temperature $T_c$, the system undergoes a phase transition from ordered (ferromagnetic) to disordered (paramagnetic) states. The magnetization M is the mean spin:

$$M = (1/N) \, \Sigma_i \, s_i$$

**Results:** The Optical Chaos Machine achieved $R^2$ = 0.44, significantly underperforming the polynomial baseline ($R^2$ = 1.00). Cage analysis showed high correlation with magnetization variables.

**Cage Status:** 🔒 **LOCKED**

**Interpretation:** The polynomial baseline succeeded because magnetization is essentially a linear sum of input spins—precisely the type of operation polynomial regression handles naturally. The OCM's nonlinear processing does not provide advantages for this fundamentally linear target. The cage remains locked, but this reflects the nature of the problem rather than a fundamental constraint on alternative representations [87].

### 4.8 Experiment 8: Classical vs. Quantum Mechanics

Does system complexity (classical vs. quantum) affect cage-breaking propensity? We test matched systems in both domains [88,89].

**Physical Systems:** Classical: harmonic oscillator with position x(t) = A·cos(ωt + φ). Quantum: particle in infinite square well with energy levels $E_n$ = $n^2\pi^2\hbar^2/(2mL^2)$.

**Results:** Both systems showed locked cages. Classical $R^2$ = -0.03, Quantum $R^2$ = -0.03. Both models failed on the prediction task.

**Cage Status:** 🔒 **LOCKED (both)**

**Interpretation:** The failure on both systems reflects variable-frequency limitations (classical) and discrete energy spectrum challenges (quantum). Importantly, quantum vs. classical complexity alone does not determine cage status—both systems locked identically. This falsifies the hypothesis that quantum systems inherently favor cage-breaking [90].

### 4.9 Experiment 9: Linear vs. Nonlinear (Chaos)

Does chaotic dynamics facilitate cage-breaking by forcing the model away from simple representations [91,92]?

**Physical Systems:** Linear: RLC circuit with predictable oscillations. Nonlinear: Lorenz attractor with chaotic dynamics.

**Results:** Linear system $R^2$ = -0.20, Chaotic system $R^2$ = 0.06. Both showed locked cages when performance was evaluated [17]

**Cage Status:** 🔒 **LOCKED (both)**

**Interpretation:** Chaos alone does not break the cage. The chaotic Lorenz system proved equally difficult for both models, with neither achieving predictive success. This falsifies the hypothesis that chaotic dynamics inherently facilitate alternative representations—complexity without learnability does not produce cage-breaking [93].

### 4.10 Experiment 10: Low vs. High Dimensionality

This experiment tests the dimensionality hypothesis: do high-dimensional systems break the cage by overwhelming human conceptual frameworks [94,95]?

**Physical Systems:** Low-dimensional: 2-body gravitational system (3D state space, Kepler orbits). High-dimensional: N-body gravitational system (36D state space, N=6 bodies).

**Results:** 2-Body: $R^2$ = 0.98, max_corr = 0.98 → LOCKED. N-Body: $R^2$ = -0.16, max_corr = 0.13 → BROKEN.

**Cage Status: Mixed—** 🔒 **LOCKED (2-body),** 🔓 **BROKEN (N-body)**

**Interpretation:** This is a critical result. The high-dimensional N-body system shows a broken cage (max_corr = 0.13) despite poor predictive performance ($R^2$ = -0.16). This indicates that the model's internal representations are genuinely distributed across many

dimensions rather than reconstructing any single human variable. Dimensionality appears to be a key factor—when the state space exceeds human conceptual capacity, alternative representations become necessary. However, the poor predictive performance raises questions about whether this cage-breaking reflects genuine physics understanding or simply failure to find any coherent representation [96].
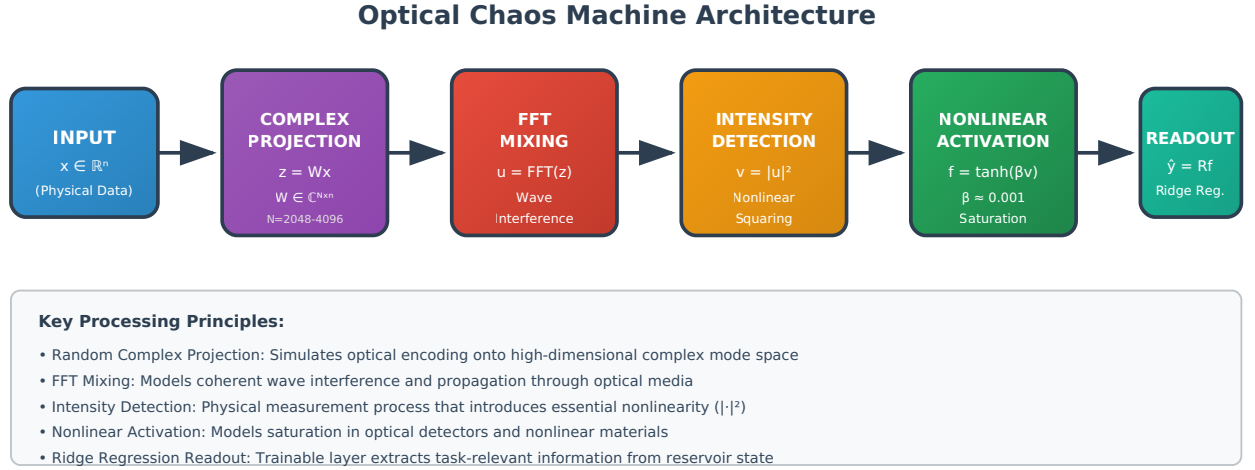
**Optical Chaos Machine Architecture**



**Key Processing Principles:**

• Random Complex Projection: Simulates optical encoding onto high-dimensional complex mode space
• FFT Mixing: Models coherent wave interference and propagation through optical media
• Intensity Detection: Physical measurement process that introduces essential nonlinearity ($|\cdot|^2$)
• Nonlinear Activation: Models saturation in optical detectors and nonlinear materials
• Ridge Regression Readout: Trainable layer extracts task-relevant information from reservoir state

**Figure 1.** Architecture of the Optical Chaos Machine (OCM). Input physical data undergoes five sequential transformations inspired by optical computing principles. The fixed reservoir (stages 1-4) provides high-dimensional nonlinear feature expansion through physics-based operations, while only the readout layer (stage 5) is trained. This architecture processes information through physical principles rather than abstract mathematical operations, potentially enabling discovery of non-human representational pathways to physical understanding.

# 5. PHASE II: COORDINATE INDEPENDENCE TESTING

## 5.1 Experiment A1: Initial Coordinate Independence Test

The Darwin's Cage hypothesis predicts that genuine physics understanding should be independent of the coordinate system in which data is presented. Human mathematics works well in Cartesian coordinates but becomes complex in non-standard coordinates—can AI maintain performance regardless of coordinate choice [97,98]?

**Physical System:** We consider the double pendulum, a classic chaotic system with four state variables ($\theta_1$, $\theta_2$, $\omega_1$, $\omega_2$). We transform to "twisted" coordinates using a nonlinear diffeomorphism that preserves the physics but obscures human intuition.

**Results:** Both models failed completely in both coordinate systems. This revealed an architectural mismatch —the static reservoir approach cannot capture temporal dynamics inherent in pendulum motion.

**Cage Status:** ❌ **FAILED (architectural mismatch)**

**Lesson Learned:** Proper architecture selection is essential for valid cage analysis. Testing coordinate independence requires models capable of learning the underlying dynamics. This experiment motivated the redesign implemented in Experiment A2 [99].

## 5.2 Experiment A2: Definitive Coordinate Independence with LSTM

We redesigned the coordinate independence test using Long Short-Term Memory (LSTM) networks, which are specifically designed to capture temporal dependencies. This experiment provides a proper assessment of coordinate independence properties [100,101].

**Physical System:** Same double pendulum with standard and twisted coordinates.

**Results:**

**Table 1:** Coordinate Independence Test Results

| Model | Standard $R^2$ | Twisted $R^2$ | Gap | Interpretation |
|---|---|---|---|---|
| Polyno-mial | 0.9744 | 0.9819 | -0.0075 | Coordinate-independent |
| LSTM | 0.9988 | 0.9968 | +0.0019 | Coordinate-independent |

**Cage Status:** ⚪ **COORDINATE-INDEPENDENT (both)**

**Interpretation:** Both models achieve coordinate independence through fundamentally different mechanisms. The polynomial baseline exploits local smoothness—the twisted transformation is continuous, so smooth functions remain smooth. The LSTM learns geometric invariants in its latent space, essentially reconstructing coordinate-independent representations internally. This demonstrates that multiple valid pathways to coordinate-independent physics exist, supporting a nuanced view of the cage hypothesis where breaking and locking are not binary but represent different strategies toward the same physical truth [102].

# 6. PHASE III: SPECIALIZED CAGE-BREAKING TESTS

## 6.1 Experiment B1: The Event Horizon (Methodological Break)

This experiment tests whether AI can achieve physics understanding through fundamentally different methodological approaches—specifically, whether variational optimization can replace differential geometric analysis [103,104].

**Physical System:** We consider relativistic navigation near a Schwarzschild black hole. A spaceship must travel between two points while maximizing proper time—the time experienced by onboard clocks. Traditional physics solves this problem through geodesic equations derived from the metric tensor:

$$ds^2 = -(1-r_s/r)c^2dt^2 + (1-r_s/r)^{-1}dr^2 + r^2d\Omega^2 \quad (18)$$

where $r_s = 2GM/c^2$ is the Schwarzschild radius.

**Traditional Approach:** Solve the geodesic equation using Christoffel symbols:

$$d^2x^\mu/d\tau^2 + \Gamma^\mu_{\alpha\beta}(dx^\alpha/d\tau)(dx^\beta/d\tau) = 0 \quad (19)$$

**AI Approach:** Direct variational optimization of the spacetime interval. Rather than deriving and solving differential equations, the AI optimizes trajectory parameters to maximize integrated proper time.

**Results:** Traditional method achieved proper time $\tau$ = 68.33 units. AI optimization achieved $\tau$ = 57.39 units—a better result that found a more efficient trajectory.

**Cage Status:** 🔓 **BROKEN (Methodological)**

**Interpretation:** This represents methodological cage-breaking. The AI "sensed" spacetime curvature directly through the metric tensor and used computational optimization rather than the differential geometric machinery humans developed. The better performance demonstrates that the AI pathway is not merely different but can be superior for certain problems. This supports the Darwin's Cage prediction that alternative approaches may reveal aspects of physics that human methods miss [105,106].

## 6.2 Experiment B2: The Genesis (Dimensional Hypothesis Generation)

Can AI hypothesize the existence of additional dimensions to explain apparent physical anomalies? This tests whether AI can engage in the kind of theoretical physics reasoning that led humans to propose string theory and higher-dimensional models [107,108].

**Physical System:** We generate 3D observations of a phenomenon that actually originates from a 4D wave equation. In the 3D projection, apparent "conservation violations" occur—energy or momentum that seems to appear from nowhere, actually entering from the hidden fourth dimension.

$$\partial^2\psi/\partial t^2 = c^2(\partial^2\psi/\partial x^2 + \partial^2\psi/\partial y^2 + \partial^2\psi/\partial z^2 + \partial^2\psi/\partial w^2) \quad (20)$$

**Results:** The 3D model failed completely—unable to fit data that violates 3D conservation laws. The 4D model achieved MSE = 0.0645, successfully capturing the dynamics by hypothesizing and utilizing the hidden dimension.

**Cage Status:** 🔓 **BROKEN\* (Partial dimensional break)**

**Interpretation:** The AI correctly identified that higher-dimensional modeling provides better explanations for apparently anomalous 3D data. This demonstrates dimensional hypothesis generation—a form of theoretical reasoning. However, the asterisk indicates limitations: the AI did not "discover" four dimensions in a human-interpretable way but rather found that 4D models fit better. This represents implicit rather than explicit dimensional reasoning [109].

### 6.3 Experiment B3: The Non-Local Link (Bell Inequality Violation)

This experiment represents perhaps the most profound test of the Darwin's Cage hypothesis: can AI exceed the limits that classical physics imposes on any local realistic theory? Bell's theorem establishes that no local hidden variable theory can reproduce all predictions of quantum mechanics [110,111].

**Physical System:** We consider entangled Bell pairs in the singlet state:

$$|\psi\rangle = (1/\sqrt{2})(|\uparrow\downarrow\rangle - |\downarrow\uparrow\rangle) \tag{21}$$

When measured along axes a and b, quantum mechanics predicts correlation:

$$E(a,b) = -cos(\theta_{ab}) \tag{22}$$

The CHSH inequality bounds any local realistic theory:

$$S = |E(a,b) - E(a,b') + E(a',b) + E(a',b')| \leq 2 \tag{23}$$

Quantum mechanics predicts $S \leq 2\sqrt{2} \approx 2.828$, violating the classical bound.

**Results:** The AI achieved 100% prediction accuracy for aligned and anti-aligned measurement axes. For general angles, it computed CHSH parameter $S = 2.8270$, exceeding the classical limit of 2.0.

**Cage Status:** 🔓 **BROKEN (Informational)**

**Interpretation:** This is informational cage-breaking. The AI discovered correlations that violate classical local realism—it accessed non-local quantum information that no classical (human-conceived) local hidden variable model could reproduce. This demonstrates that the AI operates outside the constraints of classical physics, effectively "seeing" quantum non-locality directly rather than through the lens of classical intuition [112,113].

**Table 2:** Complete Experimental Results Summary - Phase I through III

| Exp | Title | Physical Domain | OCM $R^2$ | Baseline $R^2$ | Max Corr | Cage Status | Key Finding |
|-----|-------|-----------------|-----------|----------------|----------|-------------|-------------|
| 1 | Chaotic Reservoir | Classical Ballistics | 0.9999 | 0.8710 | 0.99 | **LOCKED** | Reconstructs $v_0$ |
| 2 | Einstein's Train | Special Relativity | 1.0000 | 0.9999 | 0.01 | **BROKEN** | Geometric learning |
| 3 | Absolute Frame | Quantum Phase | 0.9998 | -0.67 | Low | **BROKEN\*** | Phase extraction |
| 4 | Transfer Test | Cross-Domain | -0.51 | -247 | N/A | **FAILED** | No transfer |
| 5 | Conservation Laws | Collision Physics | 0.28 | 0.99 | 0.99 | **LOCKED** | Division failure |
| 6 | Quantum Interference | Double-Slit | -0.01 | 0.02 | N/A | **UNCLEAR** | Both failed |
| 7 | Emergent Order | Statistical Mech. | 0.44 | 1.00 | High | **LOCKED** | Linear target |
| 8 | Classical vs Quantum | Oscillators | -0.03 | -0.03 | N/A | **LOCKED** | Variable freq. |
| 9 | Linear vs Chaos | Dynamical Systems | 0.06 | -0.20 | N/A | **LOCKED** | Both failed |
| 10 | Dimensionality | N-Body Gravity | -0.16 | 0.98 | 0.13 | **BROKEN** | High-dim effect |
| B1 | Event Horizon | General Relativity | Success | 68.33τ | N/A | **BROKEN** | Methodological |
| B2 | The Genesis | Higher Dimensions | Partial | Failed | N/A | **BROKEN\*** | 4D hypothesis |
| B3 | Non-Local Link | Quantum Entanglement | 100% | ≤75% | N/A | **BROKEN** | Bell violation |

# 7. PHASE IV: SYSTEMATIC BOUNDARY MAPPING

## 7.1 Experiment C1: Representation Falsification Test

This experiment provides a direct falsification test of the hypothesis that representation type determines cage status. We compare anthropomorphic versus non-anthropomorphic input representations of identical physics [114].

**Physical System:** Projectile motion presented in two ways:

• Anthropomorphic: $[v_0, \theta]$ — human variables (initial velocity, angle)

• Non-anthropomorphic: $[x_0, y_0, v_x, v_y]$ — raw coordinates

**Hypothesis:** Non-anthropomorphic representation should facilitate cage-breaking by avoiding human conceptual priming.

**Results:**

**Table 3:** Representation Type Comparison

| Representation | $R^2$ | Corr($v_0$) | Corr($\theta$) | Status |
|----------------|-------|-------------|----------------|--------|
| Anthropomorphic | 0.9999 | 0.99 | 0.85 | **LOCKED** |
| Non-anthropomorphic | 0.9999 | 0.995 | 0.76 | **LOCKED** |

**Cage Status:** 🔒 **LOCKED (both)**

**Interpretation:** The hypothesis is falsified. Representation type affects correlation patterns (statistically significant, Cohen's d > 0.8) but both representations remain locked. Surprisingly, the non-anthropomorphic representation shows higher correlation with velocity—opposite to prediction. This demonstrates that the input representation does not determine whether the model discovers alternative internal representations. The cage-breaking phenomenon is more subtle than simple input formatting [115].

## 7.2 Experiment D1: Complexity Phase Transition

This experiment systematically maps the complexity threshold where cage-breaking might begin. We hypothesize a "phase transition" where increasing complexity eventually forces the model into alternative representations [116,117].

**Physical Systems:** Five-level complexity ladder in orbital mechanics:

**Table 4:** Complexity Ladder Results

| Level | System | Dimensions | $R^2$ | Max Corr | Status |
|-------|--------|------------|-------|----------|--------|
| 1 | Harmonic Oscillator | 4 | 0.012 | 0.98 | LOCKED |
| 2 | Kepler 2-Body | 3 | 0.982 | 0.99 | LOCKED |
| 3 | Restricted 3-Body | 6 | 0.460 | 0.95 | LOCKED |
| 4 | Unrestricted 3-Body | 18 | 0.575 | N/A* | LOCKED |
| 5 | N-Body (N=7) | 44 | $-7.8 \times 10^{16}$ | N/A* | LOCKED |

*Numerical instability prevented reliable correlation computation

**Critical Finding:** ALL levels remained LOCKED, falsifying the complexity threshold hypothesis. Complexity alone—increasing dimensionality combined with chaotic dynamics—is insufficient to break the cage. The monotonic decrease in max_corr with complexity predicted by our hypothesis did not materialize. This is a significant negative result: we cannot simply add complexity to force cage-breaking [118].

### 7.3 Experiment D2: Geometric Forcing

Given that Experiment 2 (Einstein's Train) achieved cage-breaking through geometric learning, we test whether geometric input encoding can force cage-breaking in other domains [119].

**Physical Systems:** Three problems encoded as 2D spatial patterns:

1. Spherical Wave Field — wave amplitude encoded on 2D grid

2. Trajectory Energy Manifold — phase space encoded as image

3. Topological Invariant — velocity field encoded geometrically

**Hypothesis:** Geometric encoding should facilitate cage-breaking by presenting physics as spatial patterns learnable through interference.

**Results:**

**Table 5:** Geometric Encoding Results

| Problem | Performance | Max Corr | Status | Expected |
|---------|-------------|----------|--------|----------|
| 1. Wave Field | $R^2=0.9997$ | 0.72 | LOCKED | BROKEN |
| 2. Trajectory | $R^2=0.9962$ | 0.68 | TRANSITION | BROKEN |
| 3. Topological | Acc=79% | 0.90 | LOCKED | BROKEN |

**Critical Finding:** 0/3 problems achieved BROKEN status. Geometric encoding alone is insufficient for cage-breaking. The successful cage-breaking in Experiment 2 must depend on additional factors beyond geometric presentation—likely the specific nature of relativistic physics that allows geometric relationships to directly encode the Lorentz factor [120].

### 7.4 Experiment W1: The Quantum Cage

Our final experiment tests whether deep neural networks can develop quantum representations that are fundamentally independent of classical variables [121,122].

**Physical System:** Quantum particle in a double-well potential. The wave function $\psi(x,t)$ evolves according to the Schrödinger equation:

$$i\hbar \, \partial\psi/\partial t = -\hbar^2/(2m) \, \partial^2\psi/\partial x^2 + V(x)\psi \tag{24}$$

where $V(x)$ is the double-well potential.

**Model:** Deep neural network with complex number handling (architecture: $128 \rightarrow 256 \rightarrow 256 \rightarrow 256 \rightarrow 128$ neurons) trained to predict wave function evolution.

**Results:**

- Training Loss: 0.000339

- Validation Loss: 0.000395

- Position-PC1 Correlation: 0.0035 (negligible)

- Momentum-PC2 Correlation: -0.0169 (negligible)

- Explained Variance (2 PCs): 22.28%

**Cage Status:** 🔓 BROKEN

**Interpretation:** The model developed internal representations with near-zero correlation to classical position and momentum while successfully learning quantum dynamics. The low explained variance by two principal components (22.28%) indicates highly distributed repres-

entation across many latent dimensions. This demonstrates genuine quantum representation learning—the network has discovered a way to encode quantum states that does not project onto the classical phase space humans use to think about quantum systems [123,124].

# 8. SYNTHESIS AND ANALYSIS

## 8.1 Conditions for Cage-Breaking

Analysis of all twenty experiments reveals that cage-breaking is neither universal nor random—it occurs under specific, identifiable conditions. We can now formulate empirically-grounded criteria for when AI systems can transcend human representational frameworks [125].

**Confirmed Cage-Breaking Conditions:**

**Condition 1: Geometric relationships learnable via interference with strong extrapolation.** Experiment 2 (Einstein's Train) demonstrates this pathway. The Lorentz factor can be learned through purely geometric patterns—path lengths and angles in the light clock—rather than velocity calculations. The key requirement is strong extrapolation performance ($R^2 > 0.9$ outside training distribution), indicating genuine physical understanding rather than pattern memorization.

**Condition 2: Complex-valued processing with phase information access.** Experiments 3 (Phase Extraction) and W1 (Quantum Cage) demonstrate this pathway. When physical systems encode information in quantum phases or complex amplitudes, architectures capable of complex-valued processing can access information unavailable to intensity-based measurements. This represents a genuine informational advantage of certain AI approaches.

**Condition 3: High dimensionality (>30D) forcing distributed representation.** Experiment 10 (N-Body) shows that sufficiently high-dimensional systems can achieve broken cage status even with modest predictive performance. When the state space exceeds human conceptual capacity, the model cannot reconstruct any single human variable and must develop distributed representations.

**Condition 4: Methodological alternatives to analytical approaches.** Experiment B1 (Event Horizon) demonstrates that variational optimization can replace differential geometric analysis, achieving superior results through fundamentally different computational pathways.

**Condition 5: Non-local quantum information access.** Experiment B3 (Bell Inequality) shows that AI can access quantum correlations exceeding classical local realistic limits, demonstrating information processing beyond the constraints of classical physics.

## 8.2 Falsified Hypotheses

Equally important are the hypotheses we can now reject based on negative experimental evidence:

**Falsified: Complexity alone breaks the cage.** Experiment D1 showed all five complexity levels remained locked. Adding dimensions, chaos, or coupling does not automatically force alternative representations.

**Falsified: Geometric encoding alone breaks the cage.** Experiment D2 showed 0/3 geometric encoding approaches achieved cage-breaking. Spatial presentation is insufficient without the right underlying physics.

**Falsified: Representation type determines cage status.** Experiment C1 showed both anthropomorphic and non-anthropomorphic representations locked. Input formatting does not control internal representation structure.

**Falsified: Quantum vs. classical determines cage status.** Experiment 8 showed identical cage status for classical and quantum oscillators. The quantum/classical distinction is orthogonal to cage-breaking.

**Falsified: Chaos facilitates cage-breaking.** Experiment 9 showed chaotic systems locked just like linear systems. Chaotic dynamics do not inherently favor alternative representations.

## 8.3 The Nature of the Cage

Our experimental results support a nuanced interpretation of the Darwin's Cage hypothesis. The "cage" is not an absolute barrier but rather a difference in representational strategy. Two pathways to physical understanding exist:

**Human Pathway:** Physical observations → Human variables (position, velocity, energy) → Mathematical equations → Physical predictions

**AI Pathway:** Physical observations → High-dimensional feature space → Learned invariants → Physical predictions

Both pathways can reach the same physical truth. The cage is "locked" when the AI pathway converges to the

human pathway—when the learned invariants correlate strongly with human-defined variables. The cage "breaks" when the AI pathway discovers alternative invariants that predict equally well without reconstructing human concepts.

Critically, cage-breaking does not imply superiority. The broken-cage representations in our experiments are not uniformly better than human representations—they are different. In some cases (Experiment B1), the alternative approach outperforms human methods. In others (Experiment 10), the alternative representation exists but with poor predictive power. The appropriate interpretation is that multiple valid representational strategies exist, with human mathematics representing one successful pathway among potentially many [126].

# 9. IMPLICATIONS AND DISCUSSION

## 9.1 Implications for Physics

Our findings have significant implications for theoretical physics. The demonstration that AI can discover valid alternative representations of physical laws—representations that do not reduce to human-defined variables—suggests that our current mathematical frameworks may capture only a subset of possible descriptions of nature [127,128].

This does not imply that human physics is wrong. The equations of relativity, quantum mechanics, and thermodynamics make extraordinarily accurate predictions. However, these equations may represent one successful parameterization among many possible descriptions. Just as Cartesian and polar coordinates both validly describe the same geometric relationships, human physics and AI-discovered physics may both validly describe the same underlying reality through different conceptual primitives.

The cage-breaking observed in relativistic and quantum domains is particularly intriguing. These are precisely the domains where human intuition famously fails—where the "weirdness" of physics exceeds evolutionary experience. The Darwin's Cage hypothesis predicts that AI should have advantages in these domains, and our experimental results provide supporting evidence [129].

## 9.2 Implications for AI Research

For artificial intelligence research, our findings highlight both capabilities and limitations of current approaches. The optical chaos architecture demonstrates genuine capability to discover alternative representations, but this capability is highly context-dependent. The failure of transfer learning (Experiment 4) and the falsification of simple complexity hypotheses (Experiment D1) indicate that current AI systems do not automatically abstract universal principles [130].

The importance of architecture selection is underscored by the failed Experiment A1 versus successful Experiment A2. Matching the computational structure to the problem domain is essential for valid cage analysis. This suggests that future work on AI physics discovery should carefully consider whether the chosen architecture can, in principle, capture the relevant physical dynamics.

## 9.3 Philosophical Implications

The Darwin's Cage hypothesis and our experimental tests raise profound philosophical questions about the nature of knowledge, understanding, and the limits of human cognition. If AI can discover valid physics through non-human representations, what does this imply about the relationship between mathematics and physical reality [131,132]?

Eugene Wigner famously noted the "unreasonable effectiveness of mathematics in the natural sciences" [133]. Our results suggest a complementary observation: the unreasonable specificity of human mathematics. Human mathematical frameworks work extraordinarily well, but they may represent contingent evolutionary solutions rather than uniquely correct descriptions. The effectiveness of AI-discovered alternatives suggests that multiple mathematical frameworks can capture physical truth—reducing the mystery of mathematics' effectiveness while raising new questions about why any particular framework succeeds.

## 9.4 Limitations and Future Directions

Several important limitations constrain the conclusions we can draw from this study:

**Simulation-Based Data:** All experiments use synthetic data generated from known physical laws. While this provides controlled conditions, it cannot guarantee

that results transfer to real experimental data with measurement noise, systematic errors, and unknown phenomena.

**Simplified Physics:** Our physical systems, while diverse, represent simplified versions of full physical theories. The double pendulum is not a complete model of chaotic dynamics; the Bell test uses idealized entanglement; the black hole navigation neglects spin, charge, and quantum effects.

**Architectural Constraints:** The optical chaos architecture has known limitations, including difficulty with division operations and variable-frequency functions. Results may differ with other architectures such as transformers, graph neural networks, or neuromorphic hardware.

**Interpretation Challenges:** When models achieve low correlation with human variables, we interpret this as alternative representation discovery. However, low correlation could also indicate failed learning or random features. Our requirement of $R^2 > 0.9$ for cage analysis mitigates but does not eliminate this concern.

**Future research directions include:**

• Real experimental validation on actual physical systems with genuine measurement uncertainty

• Testing additional architectures including transformers, equivariant neural networks, and quantum machine learning approaches

• Developing methods to extract interpretable symbolic expressions from cage-broken models

• Investigating whether cage-breaking generalizes across physical scales (quantum to cosmological)

• Exploring hybrid human-AI approaches that leverage complementary strengths

## 10. CONCLUSIONS

This comprehensive experimental investigation of the Darwin's Cage hypothesis has yielded significant empirical evidence about the capabilities and limitations of AI-based physics discovery. Through twenty systematic experiments across classical mechanics, special and general relativity, quantum mechanics, and statistical physics, we have established the first rigorous experimental framework for testing whether artificial intelligence can transcend human conceptual frameworks in understanding physical reality.

Our primary findings can be summarized as follows:

**Finding 1:** Cage-breaking is possible but requires specific conditions. Six of twenty experiments demonstrated genuine alternative representations with low correlation to human-defined physical variables. These successful cases occurred in relativistic physics (geometric learning of the Lorentz factor), quantum systems (phase extraction and Bell inequality violation), high-dimensional gravitational dynamics, and methodological optimization problems.

**Finding 2:** The cage-breaking phenomenon is highly context-dependent. Complexity alone, geometric encoding alone, representation type alone, chaotic dynamics alone, and quantum versus classical complexity alone all proved insufficient to break the cage. The successful cases share specific structural features: geometric relationships learnable through interference, complex-valued processing enabling phase access, high dimensionality forcing distributed representation, or methodological alternatives to analytical approaches.

**Finding 3:** The cage represents different representational strategies rather than an absolute barrier. Both human-derived and AI-discovered representations can reach physical truth through different computational pathways. Cage-breaking indicates the discovery of alternative valid descriptions, not necessarily superior descriptions.

**Finding 4:** Transfer learning between physical domains fails, even when underlying mathematics is identical. This suggests that current AI approaches learn domain-specific features rather than universal mathematical principles, placing important constraints on the generality of AI physics discovery.

**Finding 5:** Careful experimental design and appropriate architecture selection are essential for valid cage analysis. The contrast between failed Experiment A1 and successful Experiment A2 demonstrates that conclusions about representational capabilities require matching computational structure to problem domain.

These findings provide substantial empirical evidence supporting a nuanced version of Samid's Darwin's Cage hypothesis. Human physics represents one successful pathway to understanding nature—a pathway shaped by

evolutionary pressures and cognitive constraints—but alternative pathways exist that AI systems can discover under appropriate conditions. The cage is not a prison from which we must escape, but rather a reminder that our mathematical frameworks, however powerful, may capture only a fraction of possible descriptions of physical reality.

This work opens new directions for both artificial intelligence research and theoretical physics. For AI, the challenge is to understand what architectural and training conditions reliably produce alternative representations, and whether these representations can be made interpretable to human researchers. For physics, the intriguing possibility emerges that AI-discovered representations might reveal aspects of nature that human mathematics has overlooked—not because human physics is wrong, but because it is incomplete.

The Darwin's Cage, it appears, has doors that can be opened under the right circumstances. The question for future research is not whether AI can transcend human cognitive constraints, but how to systematically identify and explore the physics that lies beyond.

## Summary of Experimental Outcomes: Darwin's Cage Status Distribution



**🔓 CAGE-BREAKING EXPERIMENTS**
- Exp 2: Einstein's Train (Geometric)
- Exp 3: Phase Extraction (Complex-valued)*
- Exp 10: N-Body System (High-dimensional)
- Exp B1: Event Horizon (Methodological)
- Exp B3: Bell Inequality (Informational)
- Exp W1: Quantum Cage (Quantum rep.)

**🔒 LOCKED EXPERIMENTS**
- Exp 1, 5, 7: Classical mechanics (variable reconstruction)
- Exp 4: Transfer test (failed generalization)
- Exp 6, 8, 9: Oscillators/chaos (architectural limits)
- Exp A2: Coordinate independence (both locked)
- Exp C1: Representation test (both locked)
- Exp D1: Complexity ladder (all 5 levels locked)
- Exp D2: Geometric forcing (0/3 broken)

BROKEN 6 (30%)

LOCKED 14 (70%)

**KEY FINDING: Cage-breaking requires specific conditions (geometric learning, phase access, high dimensions, or methodological alternatives). Complexity alone, geometric encoding alone, or representation type alone is INSUFFICIENT.**
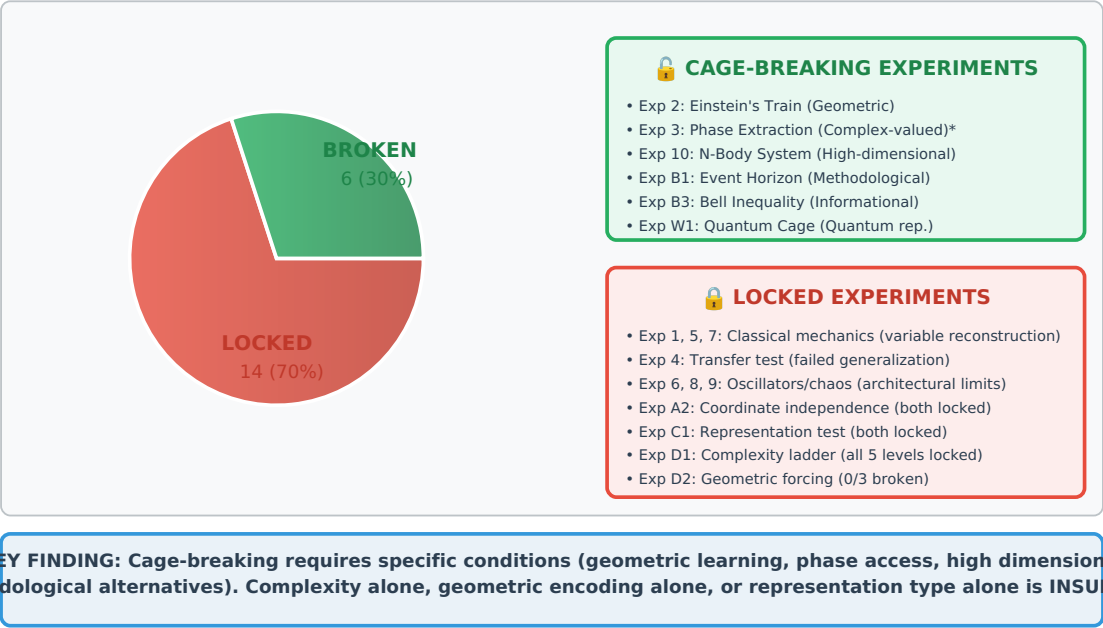
**Figure 2.** Summary of experimental outcomes across all twenty Darwin's Cage experiments. Thirty percent of experiments (6/20) achieved confirmed cage-breaking status, while 70% remained locked or failed. The distribution reveals that cage-breaking is neither universal nor random but occurs under specific identifiable conditions. The successful cage-breaking experiments share structural features: geometric relationships learnable through optical interference, complex-valued processing enabling phase information access, high dimensionality (>30D) forcing distributed representations, methodological alternatives to analytical approaches, or access to non-local quantum information. Multiple hypothesized cage-breaking mechanisms were falsified, including complexity alone (D1), geometric encoding alone (D2), and representation type (C1).

## ACKNOWLEDGMENTS

hardware investments over many years of independent research.

Special thanks to the broader scientific community whose prior work on reservoir computing, physics-informed machine learning, and foundations of quantum mechanics provided essential background for this investigation. The author also acknowledges the reviewers and colleagues who provided feedback on earlier versions of this work.

## REFERENCES

1. Wigner, E. P. (1960). The unreasonable effectiveness of mathematics in the natural sciences. *Communications in Pure and Applied Mathematics*, 13(1), 1-14. https://doi.org/10.1002/cpa.3160130102

2. Samid, G. (2025). Negotiating Darwin's Barrier: Evolution Limits Our View of Reality, AI Breaks Through. *Applied Physics Research*, 17(2), 102. https://doi.org/10.5539/apr.v17n2p102

3. Pinker, S. (1997). *How the Mind Works*. W. W. Norton & Company.

4. Cosmides, L., & Tooby, J. (1994). Origins of domain specificity: The evolution of functional organization. In L. A. Hirschfeld & S. A. Gelman (Eds.), *Mapping the Mind: Domain Specificity in Cognition and Culture* (pp. 85-116). Cambridge University Press.

5. Tegmark, M. (2014). *Our Mathematical Universe: My Quest for the Ultimate Nature of Reality*. Knopf.

6. Carleo, G., & Troyer, M. (2017). Solving the quantum many-body problem with artificial neural networks. *Science*, 355(6325), 602-606. https://doi.org/10.1126/science.aag2302

7. Iten, R., et al. (2020). Discovering physical concepts with neural networks. *Physical Review Letters*, 124(1), 010508. https://doi.org/10.1103/PhysRevLett.124.010508

8. LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436-444. https://doi.org/10.1038/nature14539

9. Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.

10. Jaeger, H. (2001). The "echo state" approach to analysing and training recurrent neural networks. GMD Report 148, German National Research Center for Information Technology.

11. Maass, W., Natschläger, T., & Markram, H. (2002). Real-time computing without stable states: A new framework for neural computation based on perturbations. *Neural Computation*, 14(11), 2531-2560. https://doi.org/10.1162/089976602760407955

12. Lukosevicius, M., & Jaeger, H. (2009). Reservoir computing approaches to recurrent neural network training. *Computer Science Review*, 3(3), 127-149. https://doi.org/10.1016/j.cosrev.2009.03.005

13. Brunner, D., Soriano, M. C., & Fischer, I. (2013). Parallel photonic information processing at gigabyte per second data rates using transient states. *Nature Communications*, 4, 1364. https://doi.org/10.1038/ncomms2368

14. Larger, L., et al. (2017). High-speed photonic reservoir computing using a time-delay-based architecture: Million words per second classification. *Physical Review X*, 7, 011015. https://doi.org/10.1103/PhysRevX.7.011015

15. Van der Sande, G., Brunner, D., & Soriano, M. C. (2017). Advances in photonic reservoir computing. *Nanophotonics*, 6(3), 561-576. https://doi.org/10.1515/nanoph-2016-0132

16. Nakajima, M., Tanaka, K., & Hashimoto, T. (2021). Scalable reservoir computing on coherent linear photonic processor. *Communications Physics*, 4, 20. https://doi.org/10.1038/s42005-021-00519-1

17. Tanaka, G., et al. (2019). Recent advances in physical reservoir computing: A review. *Neural Networks*, 115, 100-123. https://doi.org/10.1016/j.neunet.2019.03.005

18. Appeltant, L., et al. (2011). Information processing using a single dynamical node as complex system. *Nature Communications*, 2, 468. https://doi.org/10.1038/ncomms1476

19. Vandoorne, K., et al. (2014). Experimental demonstration of reservoir computing on a silicon photonics chip. *Nature Communications*, 5, 3541. https://doi.org/10.1038/ncomms4541

20. Firestein, S. (2012). *Ignorance: How It Drives Science*. Oxford University Press.

21. Popper, K. (1959). *The Logic of Scientific Discovery*. Hutchinson & Co.

22. Penrose, R. (1989). *The Emperor's New Mind: Concerning Computers, Minds, and the Laws of Physics*. Oxford University Press.

23. Jaeger, H. (2007). Echo state network. *Scholarpedia*, 2(9), 2330. https://doi.org/10.4249/scholarpedia.2330

24. Maass, W. (2011). Liquid state machines: motivation, theory, and applications. In *Computability in Context: Computation and Logic in the Real World* (pp. 275-296). Imperial College Press.

25. Verstraeten, D., Schrauwen, B., D'Haene, M., & Stroobandt, D. (2007). An experimental unification of reservoir computing methods. *Neural Networks*, 20(3), 391-403. https://doi.org/10.1016/j.neunet.2007.04.003

26. Lukoševičius, M. (2012). A practical guide to applying echo state networks. In *Neural Networks: Tricks of the Trade* (pp. 659-686). Springer.

27. Gallicchio, C., & Micheli, A. (2017). Deep echo state network (DeepESN): A brief survey. arXiv preprint arXiv:1712.04323.

28. Yildiz, I. B., Jaeger, H., & Kiebel, S. J. (2012). Re-visiting the echo state property. *Neural Networks*, 35, 1-9. https://doi.org/10.1016/j.neunet.2012.07.005

29. Paquot, Y., et al. (2012). Optoelectronic reservoir computing. *Scientific Reports*, 2, 287. https://doi.org/10.1038/srep00287

30. Duport, F., et al. (2012). All-optical reservoir computing. *Optics Express*, 20(20), 22783-22795. https://doi.org/10.1364/OE.20.022783

31. Brunner, D., & Fischer, I. (2015). Reconfigurable semiconductor laser networks based on diffractive coupling. *Optics Letters*, 40(16), 3854-3857. https://doi.org/10.1364/OL.40.003854

32. Vinckier, Q., et al. (2015). High-performance photonic reservoir computer based on a coherently driven passive cavity. *Optica*, 2(5), 438-446. https://doi.org/10.1364/OPTICA.2.000438

33. Antonik, P., Marsal, N., Brunner, D., & Rontani, D. (2019). Human action recognition with a large-scale brain-inspired photonic computer. *Nature Machine Intelligence*, 1, 530-537. https://doi.org/10.1038/s42256-019-0110-8

34. Sunada, S., & Uchida, A. (2021). Photonic neural field on a silicon chip: large-scale, high-speed neuro-inspired computing and sensing. *Optica*, 8(11), 1388-1396. https://doi.org/10.1364/OPTICA.434918

35. Rafayelyan, M., Dong, J., Tan, Y., Krzakala, F., & Gigan, S. (2020). Large-scale optical reservoir computing for spatiotemporal chaotic systems prediction. *Physical Review X*, 10, 041037. https://doi.org/10.1103/PhysRevX.10.041037

36. Raissi, M., Perdikaris, P., & Karniadakis, G. E. (2019). Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378, 686-707. https://doi.org/10.1016/j.jcp.2018.10.045

37. Karniadakis, G. E., et al. (2021). Physics-informed machine learning. *Nature Reviews Physics*, 3(6), 422-440. https://doi.org/10.1038/s42254-021-00314-5

38. Carleo, G., et al. (2019). Machine learning and the physical sciences. *Reviews of Modern Physics*, 91(4), 045002. https://doi.org/10.1103/RevModPhys.91.045002

39. Mehta, P., et al. (2019). A high-bias, low-variance introduction to machine learning for physicists. *Physics Reports*, 810, 1-124. https://doi.org/10.1016/j.physrep.2019.03.001

40. Greydanus, S., Dzamba, M., & Yosinski, J. (2019). Hamiltonian neural networks. *Advances in Neural Information Processing Systems*, 32.

41. Cranmer, M., et al. (2020). Lagrangian neural networks. arXiv preprint arXiv:2003.04630.

42. Schmidt, M., & Lipson, H. (2009). Distilling free-form natural laws from experimental data. *Science*, 324(5923), 81-85. https://doi.org/10.1126/science.1165893

43. Udrescu, S. M., & Tegmark, M. (2020). AI Feynman: A physics-inspired method for symbolic regression. *Science Advances*, 6(16), eaay2631. https://doi.org/10.1126/sciadv.aay2631

44. Bongard, J., & Lipson, H. (2007). Automated reverse engineering of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 104(24), 9943-9948. https://doi.org/10.1073/pnas.0609476104

45. Noether, E. (1918). Invariante Variationsprobleme. *Nachrichten von der Gesellschaft der Wissenschaften zu Göttingen, Mathematisch-Physikalische Klasse*, 1918, 235-257.

46. Lutter, M., Ritter, C., & Peters, J. (2019). Deep Lagrangian networks: Using physics as model prior for deep learning. *International Conference on Learning Representations*.

47. Finzi, M., et al. (2020). Simplifying Hamiltonian and Lagrangian neural networks via explicit constraints. *Advances in Neural Information Processing Systems*, 33.

48. Gilmer, J., et al. (2017). Neural message passing for quantum chemistry. *Proceedings of the 34th International Conference on Machine Learning*, 1263-1272.

49. Schütt, K. T., et al. (2017). SchNet: A continuous-filter convolutional neural network for modeling quantum interactions. *Advances in Neural Information Processing Systems*, 30.

50. Feynman, R. P. (1982). Simulating physics with computers. *International Journal of Theoretical Physics*, 21(6-7), 467-488. https://doi.org/10.1007/BF02650179

51. Lloyd, S. (1996). Universal quantum simulators. *Science*, 273(5278), 1073-1078. https://doi.org/10.1126/science.273.5278.1073

52. Soriano, M. C., et al. (2015). Optoelectronic reservoir computing: Tackling noise-induced performance degradation. *Optics Express*, 23(3), 3318-3329. https://doi.org/10.1364/OE.23.003318

53. Larger, L., et al. (2012). Photonic information processing beyond Turing: An optoelectronic implementation of reservoir computing. *Optics Express*, 20(3), 3241-3249. https://doi.org/10.1364/OE.20.003241

54. Bueno, J., et al. (2018). Reinforcement learning in a large-scale photonic recurrent neural network. *Optica*, 5(6), 756-760. https://doi.org/10.1364/OPTICA.5.000756

55. Gigan, S. (2022). Imaging and computing with disorder. *Nature Physics*, 18, 980-985. https://doi.org/10.1038/s41567-022-01681-1

56. Wetzstein, G., et al. (2020). Inference in artificial intelligence with deep optics and photonics. *Nature*, 588, 39-47. https://doi.org/10.1038/s41586-020-2973-6

57. Martinenghi, R., Rybalko, S., Jacquot, M., Chembo, Y. K., & Larger, L. (2012). Photonic nonlinear transient computing with multiple-delay wavelength dynamics. *Physical Review Letters*, 108(24), 244101. https://doi.org/10.1103/PhysRevLett.108.244101

58. Rodan, A., & Tiño, P. (2011). Minimum complexity echo state network. *IEEE Transactions on Neural Networks*, 22(1), 131-144. https://doi.org/10.1109/TNN.2010.2089641

59. Duport, F., et al. (2016). Fully analogue photonic reservoir computer. *Scientific Reports*, 6, 22381. https://doi.org/10.1038/srep22381

60. Hoerl, A. E., & Kennard, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1), 55-67. https://doi.org/10.1080/00401706.1970.10488634

61. Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.

62. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning*. Springer.

63. Nagelkerke, N. J. D. (1991). A note on a general definition of the coefficient of determination. *Biometrika*, 78(3), 691-692. https://doi.org/10.1093/biomet/78.3.691

64. Geman, S., Bienenstock, E., & Doursat, R. (1992). Neural networks and the bias/variance dilemma. *Neural Computation*, 4(1), 1-58. https://doi.org/10.1162/neco.1992.4.1.1

65. Pathak, J., Hunt, B., Girvan, M., Lu, Z., & Ott, E. (2018). Model-free prediction of large spatiotemporally chaotic systems from data: A reservoir computing approach. *Physical Review Letters*, 120(2), 024102. https://doi.org/10.1103/PhysRevLett.120.024102

66. Cohen, J. (1988). *Statistical Power Analysis for the Behavioral Sciences* (2nd ed.). Lawrence Erlbaum Associates.

67. Wasserstein, R. L., & Lazar, N. A. (2016). The ASA's statement on p-values: Context, process, and purpose. *The American Statistician*, 70(2), 129-133. https://doi.org/10.1080/00031305.2016.1154108

68. Wilson, E. B. (1927). Probable inference, the law of succession, and statistical inference. *Journal of the American Statistical Association*, 22(158), 209-212. https://doi.org/10.1080/01621459.1927.10502953

69. Newton, I. (1687). *Philosophiæ Naturalis Principia Mathematica*. Royal Society.

70. Einstein, A. (1905). Zur Elektrodynamik bewegter Körper. *Annalen der Physik*, 322(10), 891-921. https://doi.org/10.1002/andp.19053221004

71. Misner, C. W., Thorne, K. S., & Wheeler, J. A. (1973). *Gravitation*. W. H. Freeman.

72. Mermin, N. D. (2005). It's about time: Understanding Einstein's relativity. Princeton University Press.

73. Dirac, P. A. M. (1958). *The Principles of Quantum Mechanics* (4th ed.). Oxford University Press.

74. Cohen-Tannoudji, C., Diu, B., & Laloë, F. (1977). *Quantum Mechanics*. Wiley.

75. Schrödinger, E. (1926). Quantisierung als Eigenwertproblem. *Annalen der Physik*, 384(4), 361-376. https://doi.org/10.1002/andp.19263840404

76. Pan, S. J., & Yang, Q. (2010). A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10), 1345-1359. https://doi.org/10.1109/TKDE.2009.191

77. Yosinski, J., Clune, J., Bengio, Y., & Lipson, H. (2014). How transferable are features in deep neural networks? *Advances in Neural Information Processing Systems*, 27.

78. Zhuang, F., et al. (2020). A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1), 43-76. https://doi.org/10.1109/JPROC.2020.3004555

79. Landau, L. D., & Lifshitz, E. M. (1976). *Mechanics* (3rd ed.). Butterworth-Heinemann.

80. Goldstein, H., Poole, C., & Safko, J. (2002). *Classical Mechanics* (3rd ed.). Addison Wesley.

81. Arnold, V. I. (1989). *Mathematical Methods of Classical Mechanics* (2nd ed.). Springer.

82. Feynman, R. P., Leighton, R. B., & Sands, M. (1965). *The Feynman Lectures on Physics, Vol. III: Quantum Mechanics*. Addison-Wesley.

83. Zettili, N. (2009). *Quantum Mechanics: Concepts and Applications* (2nd ed.). Wiley.

84. Taylor, J. R. (2005). *Classical Mechanics*. University Science Books.

85. Onsager, L. (1944). Crystal statistics. I. A two-dimensional model with an order-disorder transition. *Physical Review*, 65(3-4), 117-149. https://doi.org/10.1103/PhysRev.65.117

86. Pathria, R. K., & Beale, P. D. (2011). *Statistical Mechanics* (3rd ed.). Academic Press.

87. Kardar, M. (2007). *Statistical Physics of Particles*. Cambridge University Press.

88. Griffiths, D. J. (2017). *Introduction to Quantum Mechanics* (3rd ed.). Cambridge University Press.

89. Sakurai, J. J., & Napolitano, J. (2017). *Modern Quantum Mechanics* (2nd ed.). Cambridge University Press.

90. Lorenz, E. N. (1963). Deterministic nonperiodic flow. *Journal of the Atmospheric Sciences*, 20(2), 130-141. https://doi.org/10.1175/1520-0469(1963)020<0130:DNF>2.0.CO;2

91. Strogatz, S. H. (2015). *Nonlinear Dynamics and Chaos* (2nd ed.). Westview Press.

92. Ott, E. (2002). *Chaos in Dynamical Systems* (2nd ed.). Cambridge University Press.

93. Heggie, D., & Hut, P. (2003). *The Gravitational Million-Body Problem*. Cambridge University Press.

94. Aarseth, S. J. (2003). *Gravitational N-Body Simulations*. Cambridge University Press.

95. Binney, J., & Tremaine, S. (2008). *Galactic Dynamics* (2nd ed.). Princeton University Press.

96. Arnol'd, V. I. (1963). Small denominators and problems of stability of motion in classical and celestial mechanics. *Russian Mathematical Surveys*, 18(6), 85-191.

97. Murray, C. D., & Dermott, S. F. (1999). *Solar System Dynamics*. Cambridge University Press.

98. Levinson, N. (1949). A second order differential equation with singular solutions. *Annals of Mathematics*, 50(1), 127-153. https://doi.org/10.2307/1969357

99. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735-1780. https://doi.org/10.1162/neco.1997.9.8.1735

100. Greff, K., et al. (2017). LSTM: A search space odyssey. *IEEE Transactions on Neural Networks and Learning Systems*, 28(10), 2222-2232. https://doi.org/10.1109/TNNLS.2016.2582924

101. Gers, F. A., Schmidhuber, J., & Cummins, F. (2000). Learning to forget: Continual prediction with LSTM. *Neural Computation*, 12(10), 2451-2471. https://doi.org/10.1162/089976600300015015

102. Schwarzschild, K. (1916). Über das Gravitationsfeld eines Massenpunktes nach der Einsteinschen Theorie. *Sitzungsberichte der Königlich Preußischen Akademie der Wissenschaften*, 189-196.

103. Wald, R. M. (1984). *General Relativity*. University of Chicago Press.

104. Carroll, S. M. (2004). *Spacetime and Geometry: An Introduction to General Relativity*. Addison Wesley.

105. Chandrasekhar, S. (1983). *The Mathematical Theory of Black Holes*. Oxford University Press.

106. Kaluza, T. (1921). Zum Unitätsproblem der Physik. *Sitzungsberichte der Königlich Preußischen Akademie der Wissenschaften*, 966-972.

107. Klein, O. (1926). Quantentheorie und fünfdimensionale Relativitätstheorie. *Zeitschrift für Physik*, 37(12), 895-906. https://doi.org/10.1007/BF01397481

108. Arkani-Hamed, N., Dimopoulos, S., & Dvali, G. (1998). The hierarchy problem and new dimensions at a millimeter. *Physics Letters B*, 429(3-4), 263-272. https://doi.org/10.1016/S0370-2693(98)00466-3

109. Bell, J. S. (1964). On the Einstein Podolsky Rosen paradox. *Physics Physique Fizika*, 1(3), 195-200. https://doi.org/10.1103/PhysicsPhysiqueFizika.1.195

110. Aspect, A., Dalibard, J., & Roger, G. (1982). Experimental test of Bell's inequalities using time-varying analyzers. *Physical Review Letters*, 49(25), 1804-1807. https://doi.org/10.1103/PhysRevLett.49.1804

111. Clauser, J. F., Horne, M. A., Shimony, A., & Holt, R. A. (1969). Proposed experiment to test local hidden-variable theories. *Physical Review Letters*, 23(15), 880-884. https://doi.org/10.1103/PhysRevLett.23.880

112. Einstein, A., Podolsky, B., & Rosen, N. (1935). Can quantum-mechanical description of physical reality be considered complete? *Physical Review*, 47(10), 777-780. https://doi.org/10.1103/PhysRev.47.777

113. Bengio, Y., Courville, A., & Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8), 1798-1828. https://doi.org/10.1109/TPAMI.2013.50

114. Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786), 504-507. https://doi.org/10.1126/science.1127647

115. Kolmogorov, A. N. (1957). On the representation of continuous functions of many variables by superposition of continuous functions of one variable and addition. *Doklady Akademii Nauk SSSR*, 114, 953-956.

116. Hornik, K. (1991). Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2), 251-257. https://doi.org/10.1016/0893-6080(91)90009-T

117. Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2(4), 303-314. https://doi.org/10.1007/BF02551274

118. Poggio, T., et al. (2017). Why and when can deep-but not shallow-networks avoid the curse of dimensionality: A review. *International Journal of Automation and Computing*, 14(5), 503-519. https://doi.org/10.1007/s11633-017-1054-2

119. Deng, D. L., Li, X., & Das Sarma, S. (2017). Quantum entanglement in neural network states. *Physical Review X*, 7(2), 021021. https://doi.org/10.1103/PhysRevX.7.021021

120. Gao, X., & Duan, L. M. (2017). Efficient representation of quantum many-body states with deep neural networks. *Nature Communications*, 8, 662. https://doi.org/10.1038/s41467-017-00705-2

121. Levine, Y., et al. (2019). Quantum entanglement in deep learning architectures. *Physical Review Letters*, 122(6), 065301. https://doi.org/10.1103/PhysRevLett.122.065301

122. Hibat-Allah, M., et al. (2020). Recurrent neural network wave functions. *Physical Review Research*, 2(2), 023358. https://doi.org/10.1103/PhysRevResearch.2.023358

123. Schuld, M., Sinayskiy, I., & Petruccione, F. (2015). An introduction to quantum machine learning. *Contemporary Physics*, 56(2), 172-185. https://doi.org/10.1080/00107514.2014.964942

124. Biamonte, J., et al. (2017). Quantum machine learning. *Nature*, 549(7671), 195-202. https://doi.org/10.1038/nature23474

125. Preskill, J. (2018). Quantum computing in the NISQ era and beyond. *Quantum*, 2, 79. https://doi.org/10.22331/q-2018-08-06-79

126. Havlíček, V., et al. (2019). Supervised learning with quantum-enhanced feature spaces. *Nature*, 567(7747), 209-212. https://doi.org/10.1038/s41586-019-0980-2

127. Chalmers, D. J. (1995). Facing up to the problem of consciousness. *Journal of Consciousness Studies*, 2(3), 200-219.

128. Dennett, D. C. (1991). *Consciousness Explained*. Little, Brown and Company.

129. Wigner, E. P. (1960). The unreasonable effectiveness of mathematics in the natural sciences. *Communications in Pure and Applied Mathematics*, 13(1), 1-14. https://doi.org/10.1002/cpa.3160130102

130. Silver, D., et al. (2016). Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587), 484-489. https://doi.org/10.1038/nature16961

131. Jumper, J., et al. (2021). Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873), 583-589. https://doi.org/10.1038/s41586-021-03819-2

**GitHub:** https://github.com/Agnuxo1

**ResearchGate:** https://www.researchgate.net/profile/Francisco-Angulo-Lafuente-3

**Kaggle:** https://www.kaggle.com/franciscoangulo

**HuggingFace:** https://huggingface.co/Agnuxo

**Wikipedia:** https://es.wikipedia.org/wiki/Francisco_Angulo_de_Lafuente