

# THEOS: Complete Mathematical Formalization

---

## Final Polished Version (Perplexity-Reviewed)

---

### Preamble

---

This document presents the complete mathematical formalization of THEOS in two parts:

- **Part A (Sections 1–6):** Core definitions, assumptions, and theorems. These establish the fundamental mathematical structure of the THEOS reasoning system.
- **Part B (Sections 7–14):** Additional constructions, sufficient conditions, and illustrative parameter settings that demonstrate the assumptions are plausible and show how they can be met in practice.

Readers should interpret Part A as the rigorous mathematical foundation and Part B as supporting evidence and practical guidance.

---

# PART A: RIGOROUS FOUNDATION

---

## Section 1: State Spaces and Spiral Dynamics

---

### Assumption 1.1: Metric Spaces and Completeness

Let  $\mathcal{O}, \mathcal{I}, \mathcal{A}, \mathcal{D}, \mathcal{F}, \mathcal{W}$  be complete metric spaces with metrics  $d_O, d_I, d_A, d_D, d_F, d_W$  respectively.

Define the combined state space:  $\mathcal{S} = \mathcal{I} \times \mathcal{A} \times \mathcal{D} \times \mathcal{F} \times \mathcal{W}$

with product metric:  $d_S((I, A, D, \Phi, \gamma), (I', A', D', \Phi', \gamma')) = d_I(I, I') + d_A(A, A') + d_D(D, D') + d_F(\Phi, \Phi') + d_W(\gamma, \gamma')$

We use  $\|\cdot\|$  for norms on  $\mathcal{D}$  and  $d_X$  for metrics on each space  $X$ .

### Definition 1.2: Phase Operators

For induction, abduction, deduction, contradiction and wisdom update, define:  $\sigma_I : \mathcal{O} \times \mathcal{F} \rightarrow \mathcal{I}$ ,  $\sigma_A : \mathcal{I} \times \mathcal{W} \rightarrow \mathcal{A}$ ,  $\sigma_D : \mathcal{A} \rightarrow \mathcal{D}$ ,  $\text{Contr} : \mathcal{D} \times \mathcal{D} \rightarrow \mathcal{F}$ ,  $\text{Upd}_\gamma : \mathcal{W} \times \mathcal{Q} \times \mathcal{D} \times \mathcal{F} \rightarrow \mathcal{W}$

Let  $q \in \mathcal{Q}$  be a fixed query and  $(O_n)_{n \geq 0}$  an observation process (possibly deterministic given  $q$ ).

### Definition 1.3: THEOS Cycle Map

Given dual deductions  $D_L, D_R$  (specified in Section 2 below), define the cycle map  $T_q : \mathcal{S} \rightarrow \mathcal{S}$  by:  $T_q(S_n) = S_{n+1} = (I_{n+1}, A_{n+1}, D_{n+1}, \Phi_{n+1}, \gamma_{n+1})$

where, for the next observation  $O'$  (or sufficient statistic of it):  $I_{n+1} = \sigma_I(O', \Phi_n)$ ,  $A_{n+1} = \sigma_A(I_{n+1}, \gamma_n)$ ,  $D_{n+1} = \sigma_D(A_{n+1})\Phi_{n+1} = \text{Contr}(D_L^{(n+1)}, D_R^{(n+1)})$ ,  $\gamma_{n+1} = \text{Upd}_\gamma(\gamma_n, q, D_n, \Phi_n)$

## Assumption 1.4: Contractive Spiral

This is an assumption to be checked or approximated in concrete implementations, not automatically guaranteed.

For the fixed query  $q$  and observation process  $(O_n)$ , there exists  $\rho \in (0, 1)$  such that:  $d_S(T_q(s), T_q(s')) \leq \rho d_S(s, s')$  for all  $s, s' \in \mathcal{S}$

## Theorem 1.5: Spiral Convergence

Under Assumption 1.4 (contractive spiral), for any initial state  $S_0 \in \mathcal{S}$ , the sequence:  $S_{n+1} = T_q(S_n)$ ,  $n \geq 0$

converges to a unique fixed point  $S^* \in \mathcal{S}$ , and:  $d_S(S_n, S^*) \leq \rho^n d_S(S_0, S^*)$

**Proof:** This is the Banach fixed-point theorem applied to the complete metric space  $(\mathcal{S}, d_S)$ .

---

## Section 2: Dual Engines and Contradiction

### Definition 2.1: Dual Engines

Let  $f_L, f_R : \mathcal{I} \times \mathcal{A} \times \mathcal{W} \rightarrow \mathcal{D}$  be the left (constructive) and right (critical) deduction functions.

At cycle  $n$ , given  $(I_n, A_n^{(L)}, A_n^{(R)}, \gamma_n)$ , define:  $D_L^{(n)} = f_L(I_n, A_n^{(L)}, \gamma_n)$ ,  $D_R^{(n)} = f_R(I_n, A_n^{(R)}, \gamma_n)$

where  $D_L^{(n)}$  and  $D_R^{(n)}$  are the left and right deductions, distinct from the combined deduction  $D_n$  in the state vector.

### Definition 2.2: Multi-Axis Contradiction

Let  $\Delta_F, \Delta_N, \Delta_C, \Delta_D : \mathcal{D} \times \mathcal{D} \rightarrow [0, \infty)$  be four discrepancy functionals measuring factual, normative, constraint, and distributional differences.

For weights  $\alpha_i \geq 0$ ,  $\sum_{i=1}^4 \alpha_i = 1$ , define the contradiction magnitude:  $\Phi_n = \Phi(D_L^{(n)}, D_R^{(n)}) := \sqrt{\alpha_1 \Delta_F^2 + \alpha_2 \Delta_N^2 + \alpha_3 \Delta_C^2 + \alpha_4 \Delta_D^2}$

## Definition 2.3: Productive Contradiction

Let  $\text{IG}(\Phi_n) \geq 0$  denote an information-gain functional (e.g., KL divergence between posteriors) and  $\text{Cost}_{\text{resolve}}(\Phi_n) > 0$  the cost of resolving the contradiction.

Define:  $Q(\Phi_n) = \frac{\text{IG}(\Phi_n)}{\text{Cost}_{\text{resolve}}(\Phi_n)}$

A contradiction is called **productive** if  $Q(\Phi_n) > \tau$  for some domain-dependent threshold  $\tau > 0$ .

## Assumption 2.4: Error Dynamics with Productive Contradiction

Let  $D^* \in \mathcal{D}$  denote the (possibly implicit) ground-truth prediction and:  $e_n = d_D(D_n, D^*)$

the prediction error at cycle  $n$ . There exist constants  $0 < \rho_1 < \rho_2 < 1$  such that:

- If  $Q(\Phi_n) \leq \tau$  (non-productive contradiction), then  $e_{n+1} \leq \rho_2 e_n$
- If  $Q(\Phi_n) > \tau$  (productive contradiction), then  $e_{n+1} \leq \rho_1 e_n$

## Assumption 2.5: Single-Engine Baseline

A single-engine baseline process satisfies:  $e_{n+1}^{\text{single}} \leq \bar{\rho} e_n^{\text{single}}$

for some  $\bar{\rho} \in (\rho_1, \rho_2)$ .

## Assumption 2.6: Frequency of Productive Contradictions

The existence of a positive lower bound  $p > 0$  on productive contradiction frequency is an empirical and design condition; Theorem 2.7 is conditional on this.

Along the dual-engine trajectory, productive contradictions occur with asymptotic frequency  $p > 0$  in the sense that:  $\liminf_{N \rightarrow \infty} \frac{1}{N} |\{0 \leq n < N : Q(\Phi_n) > \tau\}| \geq p$

## Theorem 2.7: Asymptotic Dual-Engine Advantage

Under Assumptions 2.4–2.6, the dual-engine error sequence satisfies:  $\limsup_{n \rightarrow \infty} (e_n^{\text{dual}})^{1/n} \leq \rho_{\text{eff}}$  with  $\rho_{\text{eff}} = \rho_1^p \rho_2^{1-p} < \bar{\rho}$

while the single-engine process satisfies:  $\limsup_{n \rightarrow \infty} e_n^{\text{single}} \leq \rho$

**Under these conditions**, the dual engine achieves strictly faster asymptotic geometric convergence than the single engine.

**Proof:** Over  $N$  steps, at least  $pN$  indices have contraction factor at most  $\rho_1$ , and at most  $(1 - p)N$  have at most  $\rho_2$ . Thus:  $e_N^{\text{dual}} \leq e_0^{\text{dual}} \rho_1^{pN} \rho_2^{(1-p)N} = e_0^{\text{dual}} (\rho_1^p \rho_2^{1-p})^N$

Taking  $N$ -th roots and limsup yields the claim. The single-engine part is immediate from the assumed bound.

---

## Section 3: Abduction as Optimization and Bracket

### Definition 3.1: Abduction Quality

For fixed patterns  $I \in \mathcal{I}$ , define explanatory power  $\text{EP}(H, I) \geq 0$  and complexity  $\text{Complexity}(H) \geq 0$  (e.g., via Minimum Description Length), and let  $\lambda > 0$ .

The abduction quality is:  $Q_A(H; I) = \frac{\text{EP}(H, I)}{1 + \lambda \text{Complexity}(H)}$

### Assumption 3.2: True Hypothesis and Nested Feasible Sets

There exists a (possibly idealized) ground-truth hypothesis  $H^* \in \mathcal{H}$  and a nested sequence  $(\mathcal{H}_n)_{n \geq 0}$  with  $\mathcal{H}_n \subseteq \mathcal{H}$  such that:

1.  $H^* \in \mathcal{H}_n$  for all  $n$
2.  $\mathcal{H}_{n+1} \subseteq \mathcal{H}_n$  for all  $n$
3. The update  $\mathcal{H}_n \rightarrow \mathcal{H}_{n+1}$  is induced by new evidence/contradictions and removes only hypotheses inconsistent with that evidence

### Definition 3.3: Left/Right Abductions and Quality Bracket

On cycle  $n$ , define:  $A_n^{(L)} \in \arg \max_{H \in \mathcal{H}_n} Q_A(H; I_n)$ ,  $A_n^{(R)} \in \arg \min_{H \in \mathcal{H}_n} Q_A(H; I_n)$

The bracket  $[A_n^{(R)}, A_n^{(L)}]$  is an interval in quality space (the scalar values  $Q_A$ ), not necessarily an interval in hypothesis coordinates. Let:  $q_n^{\max} = Q_A(A_n^{(L)}; I_n)$ ,  $q_n^{\min} = Q_A(A_n^{(R)}; I_n)$ ,  $\text{Width}_n = q_n^{\max} - q_n^{\min}$

## Assumption 3.4: Monotone Tightening of the Quality Bracket

There exists  $\eta \in (0, 1)$  such that for all  $n$ :  $\$Width_{n+1} \leq \eta Width_n\$$

## Theorem 3.5: Exponential Convergence in Quality

**Under Assumptions 3.2–3.4:**  $\$Width_n \leq \eta^n Width_0 \rightarrow 0 \text{ as } n \rightarrow \infty\$$

Moreover, any accumulation point  $H_\infty$  of the maximizing (or minimizing) sequence satisfies:  $\$Q_A(H_\infty; I_\infty) = \lim_{n \rightarrow \infty} q_n^{\max} = \lim_{n \rightarrow \infty} q_n^{\min}\$$

where  $I_\infty$  is the limit pattern state if it exists and is unique.

**Proof:** Immediate from the geometric bound on  $Width_n$ . The equality of the limits for accumulation points follows from continuity of  $Q_A$  under appropriate regularity.

## Assumption 3.6: Local Regularity of $Q_A$

There exists a strictly increasing function  $g : [0, \infty) \rightarrow [0, \infty)$  with  $g(0) = 0$  and a neighborhood  $\mathcal{U}$  of  $H^*$  such that for all  $H, H' \in \mathcal{U}$ :  $\$d_A(H, H') \leq g(|Q_A(H; I_\infty) - Q_A(H'; I_\infty)|)\$$

## Corollary 3.7: Abduction Convergence in Hypothesis Space

**Under Assumptions 3.2–3.6**, if  $H^*$  is the unique maximizer of  $Q_A(\cdot; I_\infty)$  in  $\mathcal{U}$ , then:  $\$d_A(A_n^{(L)}, A_n^{(R)}) \rightarrow 0 \text{ and } A_n^{(L)}, A_n^{(R)} \rightarrow H^*\$$

whenever the sequence eventually lies in  $\mathcal{U}$ .

---

## Section 4: Wisdom and Cost

---

### Definition 4.1: Wisdom Memory and Relevant Set

At cycle  $n$ , let:  $\$W_n = \{(q_i, H_i, \text{res}_i, \text{conf}_i)\}_{i=1}^n\$$

be the wisdom memory, and  $\text{Sim} : \mathcal{Q} \times \mathcal{Q} \rightarrow [0, 1]$  a similarity metric.

For a new query  $q$ , define the relevant subset:  $\$W_{\text{rel}}(q) = \{(q_i, H_i, \text{res}_i, \text{conf}_i) : \text{Sim}(q, q_i) > \sigma\}\$$

for some threshold  $\sigma \in (0, 1)$ .

### Definition 4.2: Per-Cycle Cost

The cost of cycle  $n$  is:  $\text{Cost}_n = \text{Cost}_{\text{base}} + \text{Cost}_{\text{engines}} + \text{Cost}_{\text{contr}} - \text{Savings}_{\text{wisdom}, n}$

where  $\text{Savings}_{\text{wisdom}, n} \geq 0$  is the computational savings attributable to using  $W_{\text{rel}}$ .

### Assumption 4.3: Expected Marginal Savings

For a given domain distribution over queries, there exist constants  $c > 0$  and  $\kappa > 0$  such that:  $\mathbb{E}[\text{Savings}_{\text{wisdom}, n+1} - \text{Savings}_{\text{wisdom}, n}] \geq c e^{-\kappa n} \text{ for all } n \geq 0$

### Theorem 4.4: Upper Bound on Expected Cost

**Under Assumption 4.3,** there exist constants  $C_1, C_2 > 0$  such that:  $\mathbb{E}[\text{Cost}_n] \leq C_1 + C_2 e^{-\kappa n} \text{ for all } n \geq 0$

**Proof:** Summing the marginal savings lower bound gives a convergent geometric series; rearranging yields a geometric upper bound on the expected cost term, up to additive constants that depend on the base and engine costs.

---

## Section 5: Governor and Halting

---

### Definition 5.1: Governor State

At cycle  $n$ , the governor state is:  $G_n = (\text{cycles\_remaining}_n, \text{budget\_remaining}_n, \text{confidence}_n, \text{contradiction\_trend}_n)$

### Definition 5.2: Objective for Halting

Let  $\hat{E}(n)$  be an estimate of expected error cost if halting at cycle  $n$ , and:  $C(n) = \sum_{k=0}^{n-1} \text{Cost}_k$

the cumulative computation cost up to  $n$ . For a cost-sensitivity parameter  $\lambda > 0$ , define:  $J(n) = \hat{E}(n) + \lambda C(n)$

### Definition 5.3: Optimal Halting Cycle

An optimal halting cycle is:  $n^* \in \arg \min_{n \in \mathbb{N}} J(n)$

In practice, the governor maintains an online estimate  $\hat{J}_n$  and halts when approximate local optimality and additional constraints are satisfied, e.g.:

- **Convergence:**  $\|D_L^{(n)} - D_R^{(n)}\| < \epsilon_1$  and  $Q(\Phi_n) < \tau$
- **Diminishing returns:**  $\text{IG}(\Phi_n)/\text{IG}(\Phi_{n-1}) < \rho_{\min}$
- **Budget exhaustion:**  $\text{cycles\_remaining}_n = 0$  or  $\text{budget\_remaining}_n < \text{Cost}_n$
- **Irreducible uncertainty:**  $\text{Entropy}(A_n) < \epsilon_2$  and  $\Phi_n > \delta_{\min}$

These operational criteria approximate the minimizer  $n^*$  of  $J(n)$  under uncertainty.

---

## Section 6: Output Rule

---

When the governor halts at cycle  $N$ , define the output:  $\$Output = \begin{cases} D_N^{(L)} & \text{if } \Phi_N < \epsilon_1 \\ \text{Blend}(D_N^{(L)}, D_N^{(R)}) & \text{if } \epsilon_1 \leq \Phi_N < \epsilon_2 \\ (D_N^{(L)}, D_N^{(R)}, \Phi_N) & \text{if } \Phi_N \geq \epsilon_2 \end{cases}$

where the blending operator is, for example:  $\$Blend(D_L, D_R) = w_L D_L + w_R D_R$ ,  $w_L = \frac{1-\Phi_N/\epsilon_2}{2}$ ,  $w_R = \frac{1+\Phi_N/\epsilon_2}{2}$

---

# PART B: SUPPORTING CONSTRUCTIONS AND ILLUSTRATIVE PARAMETERS

---

## Section 7: Productive Contradiction Emergence

---

### Proposition 7.1: Productive Contradiction Emergence from Finite Hypothesis Space

This is a sufficient condition for productive contradictions in finite settings.

**Claim:** Under the nested feasible set assumption, if new evidence eliminates at least one hypothesis per cycle, then productive contradictions must occur with positive frequency.

**Formal Statement:** Let  $\mathcal{H}_n$  be the feasible set at cycle  $n$  with  $|\mathcal{H}_n| = M_n$ .

**Assumption 7.1.1:** The evidence from cycle  $n$  is informative:  $M_{n+1} < M_n$  for all  $n$  (at least one hypothesis eliminated).

**Proposition:** If  $M_0 < \infty$  (finite initial hypothesis space) and  $M_{n+1} < M_n$  for all  $n$ , then:

1. The sequence terminates in finite time:  $\exists N$  such that  $M_N = 1$  (unique hypothesis)
2. Before termination, productive contradictions occur with frequency  $p \geq \frac{1}{M_0}$

**Proof Sketch:**

- By pigeonhole principle, eliminating at least one hypothesis per cycle means we exhaust the finite space in at most  $M_0$  cycles
- Each elimination corresponds to a contradiction that distinguished between hypotheses
- Such contradictions must be productive (they provide information that eliminates options)
- Therefore, productive contradictions occur at least once per  $M_0$  cycles, giving  $p \geq 1/M_0 > 0$

**Implication:** This shows that productive contradictions are plausible in finite hypothesis spaces and provides a constructive bound.

---

## Section 8: Local Regularity Justification

---

### Theorem 8.1: Sufficient Conditions for Local Regularity

These are sufficient, not necessary, conditions; in practice, we expect approximate regularity in many smooth Bayesian or neural abduction models.

**Definition 8.1.1:**  $Q_A$  satisfies **local regularity** near  $H^*$  if there exists a strictly increasing function  $g : [0, \infty) \rightarrow [0, \infty)$  with  $g(0) = 0$  such that:  $d_A(H, H') \leq g(|Q_A(H; I_\infty) - Q_A(H'; I_\infty)|)$

**Theorem:** Local regularity holds if:

1. **Lipschitz Explanatory Power:**  $\text{EP}(H, I)$  is Lipschitz continuous in  $H$  with constant  $L_E$
2. **Bounded Complexity:**  $\text{Complexity}(H)$  is bounded away from zero in a neighborhood of  $H^*$
3. **Non-degeneracy:** At  $H^*$ , the gradient  $\nabla_H Q_A(H^*; I_\infty) \neq 0$

**Proof Sketch:**

- By Lipschitz continuity:  $|Q_A(H) - Q_A(H')| \leq L_E \cdot d_A(H, H')$  (up to complexity term)
- By non-degeneracy:  $Q_A$  is strictly increasing near  $H^*$
- Therefore:  $d_A(H, H') \leq L_E^{-1} |Q_A(H) - Q_A(H')|$  locally
- Set  $g(x) = L_E^{-1}x$  to satisfy local regularity

**Implication:** This provides concrete, verifiable conditions for when Assumption 3.6 holds.

---

## Section 9: Nested Feasible Sets Mechanism

---

### Definition 9.1: Evidence-Driven Hypothesis Elimination

**Mechanism:** At cycle  $n$ , given contradiction  $\Phi_n$  between  $D_L^{(n)}$  and  $D_R^{(n)}$ :

1. **Contradiction Analysis:** Identify which hypotheses in  $\mathcal{H}_n$  would predict  $D_L^{(n)}$  vs.  $D_R^{(n)}$ 
  - Let  $\mathcal{H}_L = \{H \in \mathcal{H}_n : H \text{ predicts } D_L^{(n)}\}$
  - Let  $\mathcal{H}_R = \{H \in \mathcal{H}_n : H \text{ predicts } D_R^{(n)}\}$
2. **Elimination Rule:** Eliminate hypotheses inconsistent with the resolved contradiction
  - If evidence favors  $D_L^{(n)}$ , eliminate  $\mathcal{H}_R \setminus \mathcal{H}_L$
  - If evidence favors  $D_R^{(n)}$ , eliminate  $\mathcal{H}_L \setminus \mathcal{H}_R$
  - If evidence supports both partially, eliminate only those inconsistent with both
3. **Update:**  $\mathcal{H}_{n+1} = \mathcal{H}_n \setminus \{\text{eliminated hypotheses}\}$

### Theorem 9.2: Nesting and Ground Truth Preservation

**Theorem:** The evidence-driven elimination mechanism maintains nesting and ensures  $H^* \in \mathcal{H}_n$  for all  $n$ .

#### Proof Sketch:

- By construction,  $\mathcal{H}_{n+1} \subseteq \mathcal{H}_n$
- Since  $H^*$  generates ground truth, it predicts the actual outcome
- Therefore,  $H^*$  is never eliminated, so  $H^* \in \mathcal{H}_{n+1}$

**Implication:** This formalizes how contradictions drive hypothesis elimination, justifying Assumption 3.2.

---

## Section 10: Wisdom Similarity Formalization

---

### Definition 10.1: Semantic Similarity for Query Transfer

**Assumption 10.1.1:** Queries and hypotheses live in a semantic embedding space with distance metric  $d_Q$ .

**Definition:** Similarity between queries is:  $\text{Sim}(q, q_i) = \exp\left(-\frac{d_Q(q, q_i)^2}{2\sigma_q^2}\right)$

where  $\sigma_q$  is a domain-dependent bandwidth parameter.

### Theorem 10.2: Wisdom Transfer Benefit

**Claim:** If  $\text{Sim}(q, q_i) > \sigma$  (above threshold), then the hypothesis  $H_i$  from query  $q_i$  provides useful information for query  $q$ .

**Formal Statement:** Let  $H_i$  be the resolved hypothesis for query  $q_i$  and  $H_q^*$  be the ground truth for query  $q$ .

Under the assumption that similar queries have similar ground truths (Lipschitz in query space):  $d_A(H_i, H_q^*) \leq L_q \cdot d_Q(q, q_i)$

where  $L_q$  is the Lipschitz constant.

**Implication:** If  $\text{Sim}(q, q_i) > \sigma$ , then  $d_Q(q, q_i) < -2\sigma_q^2 \ln(\sigma)$ , which means:  $d_A(H_i, H_q^*) < L_q \cdot \sqrt{-2\sigma_q^2 \ln(\sigma)}$

So  $H_i$  is close enough to  $H_q^*$  to provide useful initialization, reducing exploration cost.

**Implication:** This formalizes why wisdom transfer works and when it's beneficial.

---

## Section 11: Information Gain Specification

---

### Definition 11.1: Information Gain from Contradiction

**Framework:** Maintain posterior distributions over hypotheses:

- $P_n(H|\text{history}_n)$  = posterior after  $n$  cycles

**Definition:** Information gain from contradiction  $\Phi_n$ :  $\text{IG}(\Phi_n) = \text{KL}(P_{n+1}(H|\text{history}_{n+1})\|P_n(H|\text{history}_n))$

This measures how much the posterior changes due to the contradiction.

## Theorem 11.2: Information Gain Bounds Entropy Reduction

**Claim:** If  $\text{IG}(\Phi_n) > \tau_{\text{IG}}$ , then the entropy of the posterior decreases by at least  $\tau_{\text{IG}}$ :  $H(P_{n+1}) \leq H(P_n) - \tau_{\text{IG}}$

**Proof:** By properties of KL divergence and entropy.

**Implication:** This connects information gain to hypothesis space reduction, justifying why productive contradictions drive progress.

---

## Section 12: Finite-Time Error Bounds

---

### Proposition 12.1: Illustrative Finite-Time Error Bound

This is an illustrative calculation with specific parameter choices, not a universal result.

**Claim:** Under the dual-engine setup with productive contradiction frequency  $p$ , the error at cycle  $n$  satisfies:  $e_n^{\text{dual}} \leq e_0 \cdot \rho_{\text{eff}}^n + \epsilon_{\text{residual}}$

where:

- $\rho_{\text{eff}} = \rho_1^p \rho_2^{1-p}$  (effective contraction rate)
- $\epsilon_{\text{residual}}$  is the irreducible error from unresolvable contradictions

**Illustrative Example:** For a representative choice of parameters:

- If  $\rho_1 = 0.3, \rho_2 = 0.7, p = 0.5$ :  $\rho_{\text{eff}} = 0.46$
- This means error reduces by ~54% per cycle
- To reach  $e_n < 0.01 \cdot e_0$ : need  $n \geq \log(0.01)/\log(0.46) \approx 5$  cycles

**Note:** Other domains will yield different constants depending on their specific contraction factors and productive contradiction frequencies.

---

# Section 13: Energy Efficiency Decomposition

---

## Proposition 13.1: Illustrative Energy Efficiency Parameterization

This is an illustrative parameterization; other domains will yield different constants.

**Qualitative Theorem:** If (i) dual-engine cycle count is lower by factor  $r > 1$ , and (ii) per-cycle cost is reduced by factor  $s \in (0, 1)$  via wisdom, then total energy is reduced by factor  $rs$ .

**Proof:** Total energy is proportional to (number of cycles)  $\times$  (cost per cycle). If cycles are reduced by  $r$  and per-cycle cost by  $s$ , total energy is reduced by  $rs$ .

**Illustrative Parameterization:** For a representative choice of parameters:

**1. Cycle Count Advantage:** Dual engines converge faster

- Single engine:  $N_{\text{single}} = \frac{\log(e_{\text{target}}/e_0)}{\log(\bar{\rho})}$
- Dual engine:  $N_{\text{THEOS}} = \frac{\log(e_{\text{target}}/e_0)}{\log(\rho_{\text{eff}})}$
- Example ratio:  $\frac{N_{\text{THEOS}}}{N_{\text{single}}} = \frac{\log(\bar{\rho})}{\log(\rho_{\text{eff}})} \approx 0.5$  (2x fewer cycles with  $\bar{\rho} = 0.8$ ,  $\rho_{\text{eff}} = 0.46$ )

**2. Per-Cycle Cost Advantage:** Wisdom reduces per-cycle cost

- With wisdom:  $C_{\text{THEOS}} = C_{\text{base}} + C_{\text{engines}} - \text{Savings}_{\text{wisdom}}$
- Example savings:  $\text{Savings}_{\text{wisdom}} \approx 0.3 \cdot C_{\text{base}}$  after accumulation
- Example ratio:  $\frac{C_{\text{THEOS}}}{C_{\text{single}}} \approx 0.7$  (30% cost reduction)

**3. Combined Example:**  $\frac{E_{\text{THEOS}}}{E_{\text{single}}} \approx 0.5 \times 0.7 = 0.35$

**Conclusion:** For this representative parameterization, total energy is approximately 35% of single-engine cost, corresponding to ~65% energy savings. Other domains will yield different percentages depending on their specific contraction factors, wisdom accumulation rates, and problem structure.

---

## Section 14: Irreducible Uncertainty Formalization

---

### Definition 14.1: Irreducible Uncertainty

**Claim:** A contradiction  $\Phi_n$  is irreducible if:

1. **Hypothesis Indistinguishability:** Both  $D_L^{(n)}$  and  $D_R^{(n)}$  are consistent with all available evidence
  - $\mathcal{H}_L \cap \mathcal{H}_R \neq \emptyset$  (overlapping predictions)
2. **Information Saturation:** Further cycles cannot distinguish between them
  - $IG(\Phi_{n+k}) < \tau_{IG}$  for all  $k > 0$
3. **Fundamental Ambiguity:** The ground truth itself is ambiguous
  - Multiple valid hypotheses explain the data equally well

### Theorem 14.2: Irreducible Contradictions Require Output Blending

**Claim:** When  $\Phi_n$  is irreducible, the optimal output is a weighted blend:  $\$Output = w_L D_L^{(n)} + w_R D_R^{(n)}$

where weights are proportional to posterior probabilities:  $w_L = \frac{P(D_L^{(n)} | \text{data})}{\sum_j P(D_j | \text{data})} \$$

**Implication:** This formalizes when to output contradictions vs. blended predictions (see Definition 5.2).

---

## PART C: SUMMARY AND DOMAIN APPLICABILITY

---

### Section 15: Complete Framework Overview

---

The THEOS mathematical core consists of:

1. **Rigorous Foundation** (Sections 1–6): Complete metric space formalism, dual-engine dynamics, abduction optimization, wisdom accumulation, governor control, and output rules
  2. **Supporting Constructions** (Sections 7–14): Justification of key assumptions, formalization of mechanisms, and illustrative parameterizations
3. **Key Properties:**

- **Convergence:** Spiral converges to fixed point with exponential rate (Theorem 1.5)
  - **Dual-Engine Advantage:** Faster asymptotic convergence than single engine under explicit conditions (Theorem 2.7)
  - **Abduction Bracket:** Quality bracket tightens exponentially (Theorem 3.5)
  - **Wisdom Benefit:** Expected cost bounded with exponential decay (Theorem 4.4)
  - **Energy Efficiency:** Illustrative parameterization shows ~65% energy savings for representative parameters (Proposition 13.1)
- 

## Section 16: Domain Universality (Conditional)

---

THEOS defines a well-posed reasoning process in any domain where:

1. **State spaces and operators** can be instantiated (Definitions 1.2–1.3)
2. **The listed assumptions approximately hold:**
  - Contractivity of the cycle map (Assumption 1.4)
  - Productive contradictions occur with positive frequency (Assumption 2.6)
  - Local regularity of the abduction quality function (Assumption 3.6)
  - Wisdom transfer provides cost reduction (Assumption 4.3)

In domains where these conditions are met or approximately met, THEOS provides a mathematically well-defined framework for reasoning under uncertainty.

---

## Section 17: Validation and Verification Checklist

---

To fully validate THEOS, verify:

### Theoretical Validation

- Prove Theorem 1.5 for specific problem classes (e.g., linear hypothesis spaces)
- Verify Assumption 2.4 empirically on benchmark datasets
- Confirm Assumption 3.4 (bracket tightening) holds for real abduction tasks
- Validate Assumption 4.3 (marginal savings) with wisdom accumulation data

### Empirical Validation

- Measure energy consumption: THEOS vs. single-engine on standard benchmarks
- Compare error rates on test sets (healthcare, finance, ethics domains)
- Track wisdom accumulation: cost reduction over time
- Test scalability: performance on problems of increasing size

### Implementation Validation

- Implement cycle map  $T_q$  for a concrete domain
  - Verify metric space properties hold in practice
  - Test Governor halting criteria on real queries
  - Validate output blending on irreducible contradictions
- 

## Section 18: Open Questions

---

1. **Productive Contradiction Frequency:** Is the bound  $p \geq 1/M_0$  tight? Can we improve it for infinite hypothesis spaces?
2. **Local Regularity:** Are the sufficient conditions in Theorem 8.1 necessary? What if they fail?

3. **Wisdom Transfer:** How does the Gaussian kernel choice affect transfer quality?  
Are there better similarity metrics?
  4. **Irreducible Uncertainty:** Can we characterize when contradictions are fundamentally unresolvable vs. just computationally hard?
  5. **Finite-Time Bounds:** Can we tighten the  $\epsilon_{\text{residual}}$  term in Proposition 12.1?
  6. **Energy Efficiency:** How sensitive is the energy savings to parameter choices?  
What are typical ranges across domains?
- 

## REFERENCES

---

- Banach Fixed-Point Theorem: Rudin, W. (1976). *Principles of Mathematical Analysis*
  - Information Theory: Cover, T. M., & Thomas, J. A. (2006). *Elements of Information Theory*
  - Optimal Stopping: Peskir, G., & Shiryaev, A. (2006). *Optimal Stopping and Free-Boundary Problems*
  - Abduction: Peirce, C. S. (1903). *Pragmatism as a Principle and Method of Right Thinking*
  - Hypothesis Testing: Neyman, J., & Pearson, E. S. (1933). On the problem of the most efficient tests
- 

**Document Status:** Final Polished Mathematical Formalization

**Perplexity Recommendations:** All 7 applied

**Date:** February 1, 2026

**Ready for:** Publication and peer review