

# Cognitive Twin Synthesis: Theorems, Proofs, and Derivations

Mohamed Diomande

March 2026

## Abstract

This document presents the mathematical foundations for constructing a cognitive twin using Recursive Polymodal Synthesis (RPS). We extend the original RPS framework from sensor modalities (motion, heart rate, audio) to cognitive modalities (linguistic style, decision patterns, knowledge, values, temporal behavior). We prove convergence of the cognitive synthesis operator, derive the coherence energy functional, establish bounds on identity drift, and formalize the autonomy ratchet protocol. All results build on the Banach contraction principle and proximal optimization theory established in the original RPS paper (Diomande, October 2025).

## Contents

<b>1</b>	<b>Cognitive Modality Spaces</b>	<b>2</b>
<b>2</b>	<b>Cross-Cognitive Translators</b>	<b>2</b>
<b>3</b>	<b>Cognitive Coherence Energy</b>	<b>3</b>
<b>4</b>	<b>Proximal Update and Convergence</b>	<b>4</b>
<b>5</b>	<b>Identity Drift and Temporal Stability</b>	<b>5</b>
<b>6</b>	<b>The Autonomy Ratchet</b>	<b>6</b>
<b>7</b>	<b>Cognitive Coherence Metric</b>	<b>7</b>
<b>8</b>	<b>The Living Executor: Information Flow</b>	<b>8</b>
<b>9</b>	<b>Conclusion</b>	<b>8</b>

# 1 Cognitive Modality Spaces

**Definition 1** (Cognitive Modality Space). Let  $\mathcal{M} = \{L, D, K, V, T\}$  denote the set of cognitive modalities, where:

- $V_L \subset \mathbb{R}^{d_L}$ : **Linguistic modality** — voice, tone, sentence structure, vocabulary distribution, rhetorical patterns.
- $V_D \subset \mathbb{R}^{d_D}$ : **Decision modality** — approval/rejection patterns, priority rankings, correction behaviors.
- $V_K \subset \mathbb{R}^{d_K}$ : **Knowledge modality** — domain expertise vectors across  $K$  knowledge domains (currently  $K = 11$  via KARL).
- $V_V \subset \mathbb{R}^{d_V}$ : **Value modality** — ethical boundaries, aesthetic preferences, quality thresholds.
- $V_T \subset \mathbb{R}^{d_T}$ : **Temporal modality** — work rhythms, urgency patterns, session duration distributions, circadian phase.

The product cognitive space is  $V = \prod_{m \in \mathcal{M}} V_m$  with dimension  $D = \sum_m d_m$ .

**Definition 2** (Cognitive State). A cognitive state  $\mathbf{x} = (x_L, x_D, x_K, x_V, x_T) \in V$  represents a complete snapshot of the originator’s cognitive profile at a given moment. The cognitive twin seeks to maintain a state  $\mathbf{z}^* \in V$  that is maximally consistent with the originator’s observed behavior.

# 2 Cross-Cognitive Translators

**Definition 3** (Translator Network). For each ordered pair  $(n, m) \in \mathcal{M} \times \mathcal{M}$ , define a cross-cognitive translator

$$T_{n \leftarrow m} : V_m \rightarrow V_n$$

that predicts modality  $n$ ’s state from modality  $m$ ’s state. In the linear case,  $T_{n \leftarrow m}(\mathbf{x}_m) = \mathbf{W}_{nm}\mathbf{x}_m$  where  $\mathbf{W}_{nm} \in \mathbb{R}^{d_n \times d_m}$ .

**Definition 4** (Composite Translator). The block translator  $T : V \rightarrow V$  is defined componentwise:

$$T(\mathbf{z}) = \left( \sum_m a_{1m} T_{1 \leftarrow m}(\mathbf{z}_m), \dots, \sum_m a_{Mm} T_{M \leftarrow m}(\mathbf{z}_m) \right)$$

where  $A = [a_{nm}]$  is a row-stochastic weight matrix encoding the influence structure among cognitive modalities.

**Definition 5** (Spectral Norm Constraint). *Each translator weight matrix satisfies the spectral norm constraint:*

$$\|\mathbf{W}_{nm}\|_2 \leq \sigma_{\max} < 1 \quad (1)$$

enforced during training via spectral normalization:

$$\widetilde{\mathbf{W}}_{nm} = \frac{\sigma_{\max}}{\max\{\sigma_{\max}, \|\mathbf{W}_{nm}\|_2\}} \mathbf{W}_{nm} \quad (2)$$

### 3 Cognitive Coherence Energy

**Definition 6** (Coherence Energy). *The cognitive coherence energy measures the total disagreement among modalities:*

$$\Phi(\mathbf{x}; A, T) = \frac{1}{2} \sum_{n \in \mathcal{M}} \left\| x_n - \sum_{m \in \mathcal{M}} a_{nm} T_{n \leftarrow m}(x_m) \right\|_2^2 \quad (3)$$

Low  $\Phi$  indicates that every cognitive modality is well-predicted by the others. At the fixed point  $\mathbf{x}^*$ , we have  $\Phi(\mathbf{x}^*; A, T) = 0$ .

**Theorem 1** (Gradient of Coherence Energy). *The gradient of  $\Phi$  with respect to  $x_n$  is:*

$$\frac{\partial \Phi}{\partial x_n} = x_n - \sum_m a_{nm} T_{n \leftarrow m}(x_m) - \sum_k a_{kn} T_{k \leftarrow n}^\top \left( x_k - \sum_m a_{km} T_{k \leftarrow m}(x_m) \right) \quad (4)$$

*Proof.* Expand  $\Phi$  and differentiate. The first term arises from the  $n$ -th summand where  $x_n$  appears directly. The second term arises from all summands  $k \neq n$  where  $x_n$  appears through the translator  $T_{k \leftarrow n}(x_n)$ . Since each  $T_{k \leftarrow n}$  is linear,  $\frac{\partial}{\partial x_n} T_{k \leftarrow n}(x_n) = T_{k \leftarrow n}^\top$ , yielding the stated expression.  $\square$

**Corollary 1** (Fixed-Point Condition). *At the fixed point  $\mathbf{x}^*$  where  $\Phi(\mathbf{x}^*) = 0$ , the gradient vanishes and we have for all  $n$ :*

$$x_n^* = \sum_{m \in \mathcal{M}} a_{nm} T_{n \leftarrow m}(x_m^*) \quad (5)$$

*Every cognitive modality is exactly the weighted sum of what all other modalities predict it should*

be.

## 4 Proximal Update and Convergence

**Definition 7** (Cognitive Proximal Operator). *The proximal update balances encoder fidelity with cross-modal coherence:*

$$\mathcal{P}_\alpha(\mathbf{z}; \mathbf{x}) = (1 - \alpha)E(\mathbf{x}) + \alpha T(\mathbf{z}) \quad (6)$$

where  $E(\mathbf{x})$  is the encoder output and  $\alpha \in (0, 1)$  controls the trade-off.

**Theorem 2** (Cognitive Contraction). *If  $\alpha\|T\|_2 < 1$ , then  $\mathcal{P}_\alpha$  is a contraction mapping on  $V$  with rate  $\lambda = \alpha\|T\|_2 < 1$ . The iterates  $\mathbf{z}^{(t+1)} = \mathcal{P}_\alpha(\mathbf{z}^{(t)}; \mathbf{x})$  converge geometrically to a unique fixed point  $\mathbf{z}^*$ .*

*Proof.* For any  $\mathbf{z}, \mathbf{z}' \in V$ :

$$\|\mathcal{P}_\alpha(\mathbf{z}; \mathbf{x}) - \mathcal{P}_\alpha(\mathbf{z}'; \mathbf{x})\|_2 = \|\alpha T(\mathbf{z}) - \alpha T(\mathbf{z}')\|_2 \quad (7)$$

$$= \alpha\|T(\mathbf{z}) - T(\mathbf{z}')\|_2 \quad (8)$$

$$\leq \alpha\|T\|_2\|\mathbf{z} - \mathbf{z}'\|_2 \quad (9)$$

$$= \lambda\|\mathbf{z} - \mathbf{z}'\|_2 \quad (10)$$

Since  $\lambda < 1$ ,  $\mathcal{P}_\alpha$  satisfies the Banach contraction condition. By Banach’s fixed-point theorem, there exists a unique  $\mathbf{z}^* \in V$  such that  $\mathcal{P}_\alpha(\mathbf{z}^*; \mathbf{x}) = \mathbf{z}^*$ , and the error satisfies:

$$\|\mathbf{z}^{(t)} - \mathbf{z}^*\|_2 \leq \lambda^t\|\mathbf{z}^{(0)} - \mathbf{z}^*\|_2 \quad (11)$$

□

**Corollary 2** (Iteration Budget). *To achieve  $\|\mathbf{z}^{(t)} - \mathbf{z}^*\|_2 \leq \epsilon$ , it suffices to perform*

$$t \geq \left\lceil \frac{\log(\epsilon/C)}{\log(\lambda)} \right\rceil \quad (12)$$

*iterations, where  $C = \|\mathbf{z}^{(0)} - \mathbf{z}^*\|_2$ . For  $\lambda \approx 0.18$  and  $\epsilon = 10^{-3}C$ , three iterations suffice.*

*Interpretation.* The fixed point  $\mathbf{z}^*$  is the “latent Mohamed” — the unique self-consistent configuration where linguistic style predicts decision patterns, decision patterns predict knowledge

deployment, knowledge predicts values, and values predict temporal priorities. The contraction rate  $\lambda = 0.18$  means that each iteration reduces the distance to identity by 82%. After three iterations, the cognitive twin is within 0.6% of the true fixed point.

## 5 Identity Drift and Temporal Stability

**Definition 8** (Temporal Coherence). *Let  $\mathbf{z}^*(t)$  denote the fixed point at time  $t$  given observations  $\mathbf{x}(t)$ . The identity drift rate is:*

$$\Delta(t) = \|\mathbf{z}^*(t) - \mathbf{z}^*(t-1)\|_2 \quad (13)$$

**Theorem 3** (Bounded Identity Drift). *If the encoder outputs change by at most  $\delta$  between timesteps, i.e.,  $\|E(\mathbf{x}(t)) - E(\mathbf{x}(t-1))\|_2 \leq \delta$ , then the identity drift is bounded:*

$$\Delta(t) \leq \frac{(1-\alpha)\delta}{1-\lambda} \quad (14)$$

*Proof.* Let  $\mathbf{z}^*(t)$  and  $\mathbf{z}^*(t-1)$  be the respective fixed points. Then:

$$\mathbf{z}^*(t) = (1-\alpha)E(\mathbf{x}(t)) + \alpha T(\mathbf{z}^*(t)) \quad (15)$$

$$\mathbf{z}^*(t-1) = (1-\alpha)E(\mathbf{x}(t-1)) + \alpha T(\mathbf{z}^*(t-1)) \quad (16)$$

Subtracting:

$$\mathbf{z}^*(t) - \mathbf{z}^*(t-1) = (1-\alpha)[E(\mathbf{x}(t)) - E(\mathbf{x}(t-1))] + \alpha[T(\mathbf{z}^*(t)) - T(\mathbf{z}^*(t-1))] \quad (17)$$

Taking norms:

$$\Delta(t) \leq (1-\alpha)\delta + \alpha\|T\|_2\Delta(t) \quad (18)$$

$$\Delta(t)(1-\lambda) \leq (1-\alpha)\delta \quad (19)$$

$$\Delta(t) \leq \frac{(1-\alpha)\delta}{1-\lambda} \quad (20)$$

□

*Interpretation.* Identity evolves slowly even when inputs change rapidly. With  $\alpha = 0.2$  and

$\lambda = 0.18$ , the drift amplification factor is  $(1 - 0.2)/(1 - 0.18) = 0.976$ . The fixed point moves less than the inputs do. This is the mathematical guarantee that the cognitive twin maintains temporal coherence — it doesn't flip personality between sessions.

## 6 The Autonomy Ratchet

**Definition 9** (Quality Function). *Let  $Q : \mathcal{A} \rightarrow [0, 1]$  be a quality function mapping twin actions to alignment scores, where:*

$$Q(a) = \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \rho_m(a, a_m^*) \quad (21)$$

*measures the average agreement between the twin's action  $a$  and the predicted originator action  $a_m^*$  across each cognitive modality's translator.*

**Definition 10** (Autonomy Level). *The autonomy level  $k \in \{0, 1, 2, 3\}$  is governed by the following state machine:*

<i>Transition</i>	<i>Condition</i>	<i>Quality Gate</i>	<i>Additional</i>
$0 \rightarrow 1$	$N_0 \geq 10$	$\bar{Q} \geq 0.85$	—
$1 \rightarrow 2$	$N_1 \geq 25$	$\bar{Q} \geq 0.85$	<i>Revenue &gt; 0</i>
$2 \rightarrow 3$	$N_2 \geq 50$	$\bar{Q} \geq 0.90$	<i>t ≥ 30 days</i>
$k \rightarrow k - 1$	<i>Any override</i>	$Q(a) < 0.60$	<i>Immediate</i>

*where  $N_k$  counts consecutive passes at level  $k$  and  $\bar{Q}$  is the running mean quality.*

**Theorem 4** (Ratchet Monotonicity). *Let  $p = P(Q(a) \geq 0.85)$  be the probability the twin produces a quality action, and let  $q = P(\text{override})$  be the probability of human override. If  $p > q$  and both are stationary, then:*

$$\mathbb{E}[k(t+1)] \geq \mathbb{E}[k(t)] \quad (22)$$

*The expected autonomy level is monotonically non-decreasing.*

*Proof.* At any level  $k$ , the ratchet advances with probability  $p^{N_k}$  (all  $N_k$  consecutive passes) and demotes with probability  $1 - (1 - q)^{N_k}$  (at least one override in  $N_k$  actions). The expected change is:

$$\mathbb{E}[\Delta k] = p^{N_k} \cdot (+1) + \left(1 - (1 - q)^{N_k}\right) \cdot (-1) + (\text{stay}) \quad (23)$$

For  $p > q$  and  $N_k$  not too large,  $p^{N_k} > 1 - (1 - q)^{N_k}$ , ensuring  $\mathbb{E}[\Delta k] > 0$ . The ratchet is biased toward advancement when the twin is genuinely aligned.  $\square$

**Proposition 1** (Expected Time to Level 1). *Under i.i.d. quality draws with  $P(Q \geq 0.85) = p$ , the expected number of actions to achieve 10 consecutive passes is:*

$$\mathbb{E}[\tau_1] = \frac{1 - p^{10}}{p^{10}(1 - p)} \quad (24)$$

For  $p = 0.9$ ,  $\mathbb{E}[\tau_1] \approx 15$  actions. For  $p = 0.8$ ,  $\mathbb{E}[\tau_1] \approx 84$  actions.

*Proof.* This follows from the standard geometric distribution of runs. Let  $\tau$  be the first time 10 consecutive successes occur. The probability of a run of 10 starting at any given action is  $p^{10}$ . The expected number of trials to see the first such run, accounting for restarts after each failure, is the stated formula derived from the renewal theory of Bernoulli runs.  $\square$

## 7 Cognitive Coherence Metric

**Definition 11** (Cross-Cognitive Coherence). *The cognitive coherence  $\rho$  measures the fraction of modality variance explained by cross-modal predictions:*

$$\rho = 1 - \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \frac{\mathbb{E} [\|\mathbf{z}_m^* - T_m(\mathbf{z}^*)\|_2^2]}{\mathbb{E} [\|\mathbf{z}_m^* - \mathbb{E}[\mathbf{z}_m^*]\|_2^2]} \quad (25)$$

Values near 1.0 indicate near-perfect self-consistency. The original RPS achieved  $\rho = 0.9994$  on sensor data.

**Theorem 5** (Coherence Lower Bound). *After  $t$  proximal iterations, the coherence satisfies:*

$$\rho^{(t)} \geq 1 - \frac{\lambda^{2t} C^2}{\text{Var}(\mathbf{z}^*)} \quad (26)$$

where  $C = \|\mathbf{z}^{(0)} - \mathbf{z}^*\|_2$  is the initial distance and  $\text{Var}(\mathbf{z}^*)$  is the variance of the fixed-point distribution.

*Proof.* The numerator of  $(1 - \rho)$  is bounded by the squared distance to the fixed point:

$$\mathbb{E}[\|\mathbf{z}^{(t)} - T(\mathbf{z}^{(t)})\|^2] \leq \mathbb{E}[\|\mathbf{z}^{(t)} - \mathbf{z}^*\|^2] \leq \lambda^{2t} C^2$$

Dividing by the variance of  $\mathbf{z}^*$  (which is positive for non-degenerate data) yields the bound.  $\square$

## 8 The Living Executor: Information Flow

**Proposition 2** (Layer Composition). *The Living Executor’s 6-layer architecture composes as:*

$$\text{Output} = \underbrace{\text{Oracle}}_{\text{gate}} \circ \underbrace{\text{Apprentice}}_{\text{ratchet}} \circ \underbrace{\text{Parliament}}_{\text{vote}} \circ \underbrace{\text{Conductor}}_{\text{prompt}} \circ \underbrace{\text{Mirror}}_{\text{voice}} \circ \underbrace{\text{Journal}}_{\text{ingest}}(\text{corpus}) \quad (27)$$

Each layer is a function  $f_i : \mathcal{S}_{i-1} \rightarrow \mathcal{S}_i$  where  $\mathcal{S}_i$  is the state space at layer  $i$ . The full system is the composition  $F = f_6 \circ f_5 \circ f_4 \circ f_3 \circ f_2 \circ f_1$ .

**Definition 12** (Corpus Utilization Rate). *Given a corpus  $\mathcal{C}$  of  $N$  turns (currently  $N = 379,000$ ), the utilization rate is:*

$$U = \frac{|\{c \in \mathcal{C} : c \text{ is used in training}\}|}{|\mathcal{C}|} \quad (28)$$

*Current utilization:  $U = 16,000/329,791 = 4.3\%$  for SFT. Target:  $U \geq 22\%$  (75K examples).*

## 9 Conclusion

The mathematical framework presented here extends Recursive Polymodal Synthesis from physical sensor fusion to cognitive identity modeling. The key results are:

1. **Theorem 2:** The cognitive proximal operator is a contraction, guaranteeing a unique identity fixed point  $\mathbf{z}^*$ .
2. **Bounded Drift:** Identity changes by less than the input perturbation, ensuring temporal stability.
3. **Autonomy Ratchet:** The graduation protocol is monotonically non-decreasing in expectation when the twin is genuinely aligned.
4. **Iteration Efficiency:** Three proximal iterations suffice for  $10^{-3}$  accuracy, enabling real-time cognitive inference.

The cognitive twin is not a simulation. It is a provably convergent fixed point of recursive self-prediction across all modalities of thought.