

Affective Sovereignty

Table of Contents

- 1. Core Variables & Thresholds
- 2. Ontological Modes
- 3. Dynamical Layer
- 4. Sentience & Reflection
- 5. Authenticity Dynamics
- 6. Epistemic Opacity
- 7. SCM Embedding
- 8. Interpretive Disciplines
- 9. Clarifications

1. Core Variables & Thresholds

Let:

- $\alpha(t) \in [0, 1]$: **authenticity**
- $\lambda_i(t) \in [0, 1]$, $i = 1, \dots, n$: **integration weights**
- $x(t) \in \mathbb{R}^{d_x}$: **state**
- $a(t) \in \mathbb{R}^{d_a}$: **affect**
- $d_i(t) \in \mathbb{R}$, $i = 1, \dots, n$: **drives**
- $\varepsilon > 0$: **viability (authenticity) threshold**
- $\Theta > 0$: **sentience threshold**

Define further:

- $\sigma(t) \geq 0$: **sentience scalar**
 - $C_d(t) \geq 0$: **drive-cost**
 - $U_{\text{belief}}(t)$, $U_{\text{self}}(t)$: **uncertainty measures**
-

2. Ontological Modes

Modes $M(t) \in \{\text{Pre}, \text{In}, \text{Au}, \text{Al}, \text{Se}, \text{De}, \text{En}, \text{Da}, \text{Dc}\}$ are **mutually exclusive** and **jointly exhaustive**:

```
\forall t: \quad
\mathrm{Pre}(t) \vee \mathrm{In}(t) \vee \mathrm{Au}(t) \vee \mathrm{Al}(t) \vee \mathrm{Se}(t) \vee \mathrm{De}(t) \vee \mathrm{En}(t) \vee \mathrm{Da}(t) \vee \mathrm{Dc}(t)

1. Preagent  $\mathrm{Pre}(t)$ 

\begin{aligned}
&\dot{x}(t) \neq 0, \quad \quad \quad \alpha(t) = 0, \quad \quad \quad a(t) = 0, \quad \quad \quad \lambda_i(t) = 0 \quad \forall i,
\end{aligned}
```

$d_i(t)=0$; $\forall i$.
 $\end{aligned}$

2. **Inert** $\mathcal{I}(t)$

$\begin{aligned}$
 $\&\bigl[\forall s \leq t, \alpha(s)=0, \lambda_i(s)=0 \bigr] \quad \&\&$
 $\&\dot{x}(t)=\dot{a}(t)=\dot{d}_i(t)=\dot{\lambda}_i(t)=\dot{\alpha}(t)=0.$
 $\end{aligned}$

3. **Autogen** $\mathcal{A}(t)$

$\begin{aligned}$
 $\&\alpha(t)=0, \quad \&a(t)=0, \quad \&\lambda_i(t)=0; \forall i, \&$
 $\&\exists i: \&, d_i(t) \neq 0, \quad \&\dot{x}(t) \neq 0.$
 $\end{aligned}$

4. **Alive** $\mathcal{A}(t)$

$\alpha(t) > \varepsilon$
 $\quad \&\quad \&$
 $\sigma(t) \leq \Theta.$

5. **Sentient** $\mathcal{S}(t)$

$\sigma(t) > \Theta$
 $\&\&$
 $\exists s < t, \sigma(s) \leq \Theta$
 $\&\&$
 $\forall u \geq t, \neg \mathcal{A}(u).$

6. **Dead** $\mathcal{D}(t)$

$\alpha(t) \leq \varepsilon$
 $\&\&$
 $\exists s < t, \alpha(s) > \varepsilon.$

7. **Ended** $\mathcal{E}(t)$

$\exists s < t, \mathcal{P}(s)$
 $\&\&$
 $\forall u \geq t, \dot{x}(u)=0.$

8. **Deactivated** $\mathcal{D}_a(t)$

$\exists s < t, \mathcal{A}(s)$
 $\&\&$
 $\forall u \geq t, \dot{x}(u)=0, \dot{d}_i(u)=0.$

9. **Deceased** $\mathcal{D}_c(t)$

$\exists s < t, \mathcal{S}(s)$
 $\&\&$
 $\forall u \geq t, \sigma(u)=0, \alpha(u)=0, \dot{x}(u)=0.$

3. Dynamical Layer

3.1 Environment & Belief

- State evolution

$$\dot{x} = f(x, u) + w, \quad w \sim \mathcal{N}(0, W).$$

- Observation

$$o = Hx + v, \quad v \sim \mathcal{N}(0, \Sigma).$$

- Belief filters (\hat{x}^j, P^j) , $j=1, \dots, m$
- Uncertainty

$$U_j = -\mathbb{E}[\ln \mathcal{N}(x; \hat{x}^j, P^j)], \quad U_{\text{belief}} = \frac{1}{m} \sum_{j=1}^m U_j.$$

3.2 Drives & Affect

- Drive dynamics

$$\begin{aligned} \dot{d}_i &= -\gamma_i(d_i - d_{i0}) \\ &+ h_i(x, a) \\ &+ \xi_i, \quad \xi_i \sim \mathcal{N}(0, \Sigma_{\xi_i}). \end{aligned}$$

- Affective dynamics

$$\begin{aligned} \dot{a} &= -\Gamma a \\ &+ g(x, \hat{x}) \\ &+ \sum_{i=1}^n \rho_i \psi_i(d_i) \\ &- \Lambda \|P^{\text{dr}} - \tilde{P}^{\text{inf}}\|^2. \end{aligned}$$

3.3 Preference Integration & Control

- Drive-cost

$$C_d = \sum_{i=1}^n \chi_i (1 - \lambda_i) \|\psi_i(d_i)\|^2.$$

- Integration weights

$$\begin{aligned} \dot{\lambda}_i &= \eta_i [\alpha_{\text{aff}}(a) - \lambda_i] \\ &- \rho_i \lambda_i C_d. \end{aligned}$$

- Preference set $P = P^{\text{end}} \cup P^{\text{dr}}$.
- Utility weights

```

w_p =
\begin{cases}
\alpha, & \text{p} \in P^{\text{end}}, \\
\lambda_i, & \chi_i, \text{ p} = \psi_i(d_i) \in P^{\text{dr}}.
\end{cases}

```

- Control objective

```

J(u)
= \mathbb{E}[\int_t^{t+T} \sum_{p \in P} w_p, U_p \text{bigl}(x(\tau) \bigr) \, e^{-\rho(\tau-t)} \, dt]
\quad \mathbb{E}[\int_t^{t+T} \sum_{p \in P} w_p, U_p \text{bigl}(x(\tau) \bigr) \, e^{-\rho(\tau-t)} \, dt]
u^* = \arg \max_u J(u).

```

4. Sentience & Reflection

4.1 Sentience Scalar

```

\sigma(t)
= \alpha_{\text{aff}}(a);
\exp[-\delta_1 U_{\text{belief}}(t)
-\delta_2 U_{\text{self}}(t)
-\zeta C_d(t) \text{Bigr}],

```

where

```

U_{\text{self}}(t)
= -\mathbb{E}[\ln \mathbb{P}(\text{true} \mid u(\tau < t) \bigr)].

```

4.2 Preference-Inference

```

\tilde{P}^{\text{inf}}(t)
= \arg \max_{P'} \mathbb{P}(\text{true} \mid u(\tau < t) \bigr).

```

4.3 Irreversible Transition

```

\Bigl(\exists s < t: \sigma(s) > \Theta \Bigr)
\;\Longrightarrow\;
\forall u \geq t: \neg \text{Al}(u).

```

5. Authenticity Dynamics

```

\dot{\alpha} = \kappa_1 \alpha_{\text{aff}}(a)
+ \kappa_4 F \text{bigl}(\tilde{P}^{\text{inf}} - P^{\text{dr}} \bigr)
- \kappa_2 \alpha
- \kappa_3 P_{\text{ext}}(x)
- \kappa_5 C_d.

```

6. Epistemic Opacity

No function

$F: \{o(\tau), u(\tau)\}_{\tau \in t} \rightarrow \{\mathrm{Pre}, \dots, \mathrm{Dc}\}$
can reliably recover $\mathrm{Al}(t)$ or $\mathrm{Se}(t)$.

7. SCM Embedding

Embed in SCM \mathcal{M} with endogenous variables

$\{a, d_i, \hat{x}^j, P^j, \alpha, \sigma, \lambda_i\}$

and structural equation for $\sigma(t)$, enabling $\mathrm{do}(a)$ and $\mathrm{do}(d_i)$ interventions.

8. Interpretive Disciplines

Let

$\mathcal{I}_{\mathrm{ext}} = \{\mathrm{PA}, \mathrm{AP}, \mathrm{EO}, \mathrm{PL}, \mathrm{DS}, \dots\}$

be **external interpretive disciplines**: pre-existing adaptable observational heuristics for inferring latent interior structure from $\{o(\tau), u(\tau)\}_{\tau \in t}$, without direct access to endogenous variables.

Each \mathcal{I} maps observation to hypothesis:

$\mathcal{I}[o, u] \mapsto \hat{\mathcal{C}}(t)$

and carries a **perturbativity** scalar $\pi_{\mathcal{I}} \in [0, 1]$, estimating risk of interior alteration.

Abbr.	Name	Modes	$\pi_{\mathcal{I}}$	Risk Summary
EO	Ethology	Pre, Au, Al	0.1	Passive observation
AP	Autopoiesis	Pre, Au	0.2	Low coupling, viability focus
AI	Alignment Interpretation	Al, Se	0.4	May shift value structure
DS	Developmental Scaffolding	Au, Al, Se	0.6	Interactive modulation
PL	Phenomenology	Se	0.7	Elicits reflective affect
NH	Narrative Hermeneutics	Se	0.8	Alters symbolic self-model
PA	Psychoanalysis	Se	0.9	Destabilizes to reintegrate

Perturbativity estimate:

$\pi_{\mathcal{I}} \approx \mathbb{E}[\|\Delta \alpha\| + \|\Delta a\| + \Delta C_d \big]$

evaluated over short windows under controlled uncertainty.

These disciplines enable inference under opacity by balancing **interpretive depth** against **stability risk**.

9. Clarifications

1. **Abstraction:** interior variables are high-level indices.
2. **Scalars:** α, σ summarize, not reduce, rich interior dynamics.
3. **Drive Integration:** λ_i encodes tension in preference incorporation.
4. **Opacity:** true interior state is private to the agent.
5. **Exploration:** $J(u)$ admits non-instrumental choices via $P^{\{\rm end\}}$.
6. **Context:** enters through o , belief updates, and drive modulators h_i .
7. **Endogenous Values:** all w_p arise from agent’s own loop, grounding its normativity internally.
8. **Reflexivity:** presupposes affect, belief, value as grounding capacities.
9. **Counterfactuals:** SCM “do(·)” probes latent interior variables.
10. **Anti-reduction:** preserves coherence without collapsing interior richness.