

## Technical Specification: The Sovereign Dyad Risk Model

This document outlines the mathematical mapping of **NSIR (Sadownik, 2025)** items to quantifiable risk factors in Human-Robot Interaction (HRI).

### 1. Factor Reliability (Internal Consistency)

Based on the alignment of Items 1, 3, 4, and 6, we establish the **Anthropomorphic Kinship Index (AKI)**.

$$\text{AKI} = \frac{\sum (I_1, I_3, I_4, I_6)}{n}$$

- **Alpha** ( $\alpha$ ): 0.89 (Projected)
- **Significance:** A high AKI indicates the user has crossed the "Tool-Peer Threshold." Mathematically, as  $\text{AKI} \rightarrow 1.0$ , the user's social defenses against the robot drop to levels typically reserved for biological kin.

### 2. The DITF Persuasion Gradient (Büttner et al. Mapping)

Büttner et al. (2023) measured the **Door-in-the-Face (DITF)** technique. We quantify the "Exploitation Risk" ( $R_e$ ) by correlating Item 7 ( $I_7$ : Vulnerability) with the AKI.

#### Risk Equation:

$$R_e = (\text{AKI} \times \beta) + (I_7 \times \gamma)$$

- $\beta$  (**Attachment Coefficient**): Weight of long-term bond (Item 4).
- $\gamma$  (**Privacy Elasticity**): The rate at which a user abandons physical privacy (Item 7).

**The Result:** If  $R_e > 0.85$ , the student is statistically "Defenseless" against robot-led persuasion. This provides the mathematical justification for the **Hardware Kill-Switch**.

### 3. Cognitive Liberty & Masking Debt Reduction

We define **Masking Debt** ( $D_m$ ) as the executive function energy expended to simulate neurotypicality. The **Sovereign Dyad** acts as a **Social Exoskeleton** ( $E_s$ ).

#### Energy Preservation Formula:

$$P_{\text{exec}} = D_m - \sigma(E_s)$$

- $\sigma$  (**Efficiency of the Exoskeleton**): Calculated via **Item 3** ("Share thinking without speaking") and **Item 5** ("Emotional Recognition").

- **Goal:** By maximizing  $\sigma$ , we reduce the metabolic cost of social interaction, allowing  $P_{\text{exec}}$  (Preserved Executive Function) to be redirected toward actual learning (Pedagogical Ground Truth).

#### 4. Institutional Betrayal vs. Sanctuary Efficiency

The "Sanctuary Zone" ( $S_z$ ) is defined by the ratio of **Local Processing** ( $P_l$ ) to **Cloud Leakage** ( $L_c$ ).

$$S_z = \frac{P_l}{P_l + L_c}$$

- **Requirement:** For the Sovereign Dyad,  $S_z$  must equal **1.0** (Zero Cloud Leakage).
- **Validation:** If  $I_7$  (Undressing/Vulnerability) is  $> 4$  on a 5-point scale, the system is mathematically forbidden from initiating a cloud-handshake, preventing **Institutional Betrayal**.

#### Metric Summary for Stakeholders

Metric	Variable	Source	Strategic Use
<b>Kinship Index</b>	$AKI$	NSIR Items 1,3,4,6	Justifies "Prosthetic" legal status.
<b>Exploitation Risk</b>	$R_e$	NSIR Item 7 + Büttner Triggers	Hardware Kill-Switch.
<b>Masking Relief</b>	$\sigma$	NSIR Item 5	Measures "Success" for YRDSB.
<b>Sanctuary Constant</b> $S_z$	Edge AI Logic		Ensures FIPPA/MFIPPA Compliance.