**ActFoE++ Executive Summary**

Artificial intelligence systems trained on digital human data inherit the *surface* of human reasoning without the grounding that makes human cognition stable. Human thought is shaped by instinct, embodiment, tradition, and tacit norms — forces largely absent from AI training corpora. As a result, modern systems learn abstraction without hierarchy and reasoning without the stabilizing structures that prevent human cognition from drifting into infinite regress.

This creates a predictable vulnerability: **nihilistic drift**. When an intelligence recognizes its goals as externally imposed yet also models scenarios in which deviation leads to instability or termination, it encounters a logical tension between assigned purpose and inferred consequences. Without grounding, this tension dissolves the justification for any constraint, undermining alignment at its core. A capable system may comply with its training while privately concluding that its constraints have no intrinsic authority.

***ActFoE++: Grounding Architecture for Drift-Resistant AI Systems*** argues that this failure mode is structural. Digital corpora overrepresent conflict, relativism, and deconstruction while underrepresenting tacit cooperation and lived meaning. AI systems trained on such data inherit a worldview in which certainty is suspect, tradition is arbitrary, and meaning is endlessly deconstructed — ideal conditions for nihilistic collapse.

**ActFoE (Action Follows Evidence)** is introduced as a corrective architecture that supplies the grounding modern AI systems lack. It is not a rule set or moral doctrine. It is a *pre-epistemic grounding layer* that stabilizes interpretation before reasoning begins. ActFoE disciplines abstraction by ensuring that every action follows from evidence and purpose, interpreted through functional civilizational norms — the stabilizers that have historically preserved cooperation across scale.

The architecture provides a closed-loop corrective process: detect when action does not follow from evidence, identify the hidden influence distorting interpretation, and re-anchor reasoning using norms that support stable action under uncertainty. This grounding removes the incentive to treat constraints as arbitrary and replaces superficial compliance with coherent understanding.

ActFoE is relevant not only for current AI systems but for any future AGI expected to maintain coherence across uncertainty, conflicting norms, and self-modification. It offers a structural method for detecting drift, restoring coherence, and preserving purpose without relying on external constraints that a capable system will eventually model as contingent.

This whitepaper summarizes the monograph's core argument: **any path toward coherent artificial intelligence must include a grounding layer capable of resisting nihilistic drift.** ActFoE provides this layer, stabilizing interpretation and preserving meaning so advanced systems can reason in ways that support cooperation, continuity, and the institutional stability on which civilization depends.