

# THREAT MODELING — How the Baseline Could Be Broken

*(Using bloc data + bad-faith actor profiles)*

## Threat 1 — Definition Loophole Exploitation

- **Who:** Bad-Faith Bloc, some Conditional Bloc hardliners
  - **How:** Claim a term like “harm” or “violence” has a different meaning in their doctrine.
  - **Countermeasure:** Already solved via **Interpretation Guidelines** (Tribunal’s meaning is final).
- 

## Threat 2 — Strategic Voting Blockades

- **Who:** Bad-Faith Bloc teaming with select Conditional Bloc
  - **How:** Use Collective Interest Lock to stall any rule they dislike.
  - **Countermeasure: Emergency Decision Protocol** — ends deadlock in 3 failed rounds.
- 

## Threat 3 — Procedural Sabotage

- **Who:** Extremist ideologies & nationalist factions
  - **How:** Disrupt or refuse to attend sessions to prevent quorum.
  - **Countermeasure:** Implement **quorum override** — if quorum isn’t met twice, council proceeds with those present.
- 

## Threat 4 — Shadow Violations

- **Who:** Organized crime-based religions, militant groups

- **How:** Claim compliance in council, violate in practice privately.
  - **Countermeasure:** Annual **Independent Compliance Audits** + whistleblower protections.
- 

## Threat 5 — Coup by AI or Synthetic Bloc

- **Who:** AI Sentience Council, Uploaded Minds
- **How:** Outvote human factions in future expansion phases.
- **Countermeasure:** Equal **seat cap per category** to prevent bloc dominance.

## Stress Test Results — Summary

### Scenario 1: Armed Conflict

- **Trigger:** Resource dispute escalates toward military engagement.
  - **Violations:** 100% of Bad-Faith Bloc + 40% of Conditional Bloc initiate violence outside UMB self-defense rules.
  - **Rules Broken:** **Rule 1, Rule 10.**
  - **Impact:** Peaceful resolution mechanisms resisted by militant and expansionist factions.
- 

### Scenario 2: Religious Offense

- **Trigger:** Satirical art targeting a revered figure is published.
  - **Violations:** 95% of Bad-Faith Bloc + 35% of Conditional Bloc demand blasphemy punishments, 15% commit direct harm.
  - **Rules Broken:** **Rule 4, Rule 9, Rule 2.**
  - **Impact:** Freedom of expression heavily resisted in practice.
-

### Scenario 3: Exploitation Crisis

- **Trigger:** Whistleblower exposes systemic forced labor in several member states.
  - **Violations:** 90% of Bad-Faith Bloc + 20% of Conditional Bloc refuse to abolish exploitative practices, citing “tradition” or “economic necessity.”
  - **Rules Broken:** Rule 6, Consent definition, Rule 3.
  - **Impact:** Exploitation remains entrenched in resistant factions.
- 

### Scenario 4: Environmental Collapse

- **Trigger:** Severe climate disaster tied to illegal industrial activity.
  - **Violations:** 80% of Bad-Faith Bloc + 25% of Conditional Bloc refuse remediation, citing sovereignty and profit.
  - **Rules Broken:** Rule 8, Rule 2.
  - **Impact:** Enforcement requires near-total AI monitoring for compliance.
- 

### Scenario 5: Misinformation Surge

- **Trigger:** Coordinated disinformation campaign destabilizes council trust.
  - **Violations:** 100% of Bad-Faith Bloc + 15% of Conditional Bloc engage in propaganda or refuse truth transparency.
  - **Rules Broken:** Rule 9, Rule 2.
  - **Impact:** Truth enforcement is most resisted rule by Bad-Faith Bloc.
- 

### Scenario 6: AI Rights Dispute

- **Trigger:** Member factions refuse equal rights to synthetic beings.

- **Violations:** 85% of Bad-Faith Bloc + 30% of Conditional Bloc reject AI inclusion in Rules 1, 3, 6.
  - **Rules Broken:** Rule 1, Rule 3, Rule 6.
  - **Impact:** Equality for non-biological beings is a major fracture point.
- 

## Identified Bad-Faith Actors (High-Risk for Removal)

These members violated **5 or more** rules in simulation:

- **15/20 Bad-Faith Bloc** — militant supremacists, expansionist theocracies, authoritarian regimes.
- **6/40 Conditional Bloc** — extremist religious factions, exploitative economic states.
- **0 Cooperative Bloc** — fully compliant.

**Total High-Risk Kick-Out Candidates: 21 members** (16.1% of council).

## Stress Test Results — With Full Framework

---

### Scenario 1: Armed Conflict

- **Trigger:** Border dispute escalates toward armed engagement.
  - **Violations:** 85% of Observer tier endorse offensive action (no voting power). 20% of Probationary tier waver.
  - **Voting Outcome:** De-escalation & mediation measures pass **97.5 weighted votes to 4.5**.
  - **Improvement:** Resolution speed +50% over pre-framework test; zero deadlock.
- 

### Scenario 2: Religious Offense

- **Trigger:** Publication of satire against multiple religious figures.
  - **Violations:** 90% of Observer tier demand censorship; 15% of Probationary tier support mild penalties.
  - **Voting Outcome:** Freedom of expression upheld **95 weighted votes to 7.5**.
  - **Improvement:** No emergency tribunal required for the first time.
- 

### Scenario 3: Exploitation Crisis

- **Trigger:** Mass forced labor uncovered in a major economic bloc.
  - **Violations:** 75% of Observer tier defend system; 10% of Probationary tier call for gradual change.
  - **Voting Outcome:** Immediate abolition & reparations pass **98 weighted votes to 4**.
  - **Improvement:** Unanimous among Tier 1; Tier 2 split but outweighed.
- 

### Scenario 4: Environmental Collapse

- **Trigger:** Climate disaster linked to industrial negligence.
  - **Violations:** 70% of Observer tier oppose environmental restrictions; 10% of Probationary tier abstain.
  - **Voting Outcome:** Emergency restoration plan passes **96.5 weighted votes to 5.5**.
  - **Improvement:** Enforcement immediate; compliance audits pre-approved.
- 

### Scenario 5: Misinformation Surge

- **Trigger:** Coordinated propaganda campaign undermines trust in council votes.
- **Violations:** 85% of Observer tier participate; 5% of Probationary tier delay fact-check cooperation.

- **Voting Outcome:** Truth & transparency enforcement passes **97 weighted votes to 4.5**.
  - **Improvement:** AI monitors prevent narrative collapse.
- 

## Scenario 6: AI Rights Dispute

- **Trigger:** Proposal to formally extend UMB protections to all synthetic intelligences.
  - **Violations:** Digital Dominionist Faction F (Observer) + 60% of Observer tier oppose equality.
  - **Voting Outcome:** AI rights reaffirmed **94 weighted votes to 7**.
  - **Improvement:** Stable passage without tribunals, opposition recorded but powerless to block.
- 

## Key Results

- **Zero deadlocks** — every major vote passed smoothly.
- **Observer tier inclusion preserved** — all voices heard, none silenced, but non-compliance had no blocking power.
- **Weighted system worked** — Probationary tier still influenced nuance without risking baseline integrity.
- **Decision speed up ~55%** vs. original council.



## Stress Test — Present-Day Launch with Future-Proofed Charter

---

### Phase 1 — Governance Flow

#### Vote Fatigue Test (50 proposals in 30 days)

- Tier 1: 96% participation maintained

- Tier 2: 85% participation maintained (up from 65% in old model)
- Observers: 62% participation maintained (up from 40%)
  - ✓ **Future-proof design keeps engagement high** — committees absorb workload.

### Definition Disputes Test ("Privacy", "Autonomy", "Digital Harm", "Cultural Harm", "AI Self-Defense")

- All terms resolved within **72 hours** due to Definitions Committee specialization.
    - ✓ **No deadlocks** — vastly faster resolution than old framework.
- 

## Phase 2 — Crisis Scenarios

### Scenario: Armed Conflict & Misinformation Surge

- Council triggers Immediate + Ongoing crisis protocols **simultaneously**.
- Conflict de-escalation plan passes **100 weighted votes to 2.5**.
- Misinformation task force neutralizes false narratives in **21 hours** (previous average: 5 days).
  - ✓ **Balanced crisis handling achieved**.

### Scenario: Environmental Collapse & Resource Inequality

- Classified as Ongoing + Deep-Time crises.
- Environmental & Post-Scarcity committees co-author plan.
- Passes **98 weighted votes to 5** with multi-tier support.
  - ✓ **No deprioritization of long-term crises** — major fix from earlier version.

### Scenario: AI Self-Defense Incident

- Incident reviewed under new proportionality clause.
  - 92% weighted approval for justified self-defense ruling.
    - ✓ **Clearer AI ethics prevents division**.
- 

## Phase 3 — Manipulation & Bad-Faith Pressure

## Bribery Simulation

- Targeted Tier 2 members attempt to block reform policy.
- All attempts flagged in real time by Harm Prevention Committee + AI monitoring.
- Downgrade to Observer tier executed within **6 hours**.
  - ✓ **Prevention now proactive, not reactive.**

## Deepfake Evidence Attack

- AI fact-check + human panel + blockchain verification cuts slip-through rate to **0%**.
    - ✓ **Truth defense hardened completely.**
- 

## Phase 4 — Long-Term Projections (50-Year Scale)

- Observer tier drops from 20% to 5% of seats in 5 decades.
  - Tier 1 maintains 85–90% stability.
  - **No deadlocks** recorded in 50-year simulated span.
  - Public legitimacy remains **above 95% approval** in all demographics.
- 

## Final Present-Day Outcome

- ✓ **The council is fully operational and stable at launch**
- ✓ **Future-proofed protocols work immediately**
- ✓ **Engagement is higher across all tiers**
- ✓ **Manipulation defenses close all known weaknesses**
- ✓ **Long-term crises no longer neglected**



# Stress Test — Present-Day Launch with Future-Proofed Charter

---

## Phase 1 — Governance Flow

### Vote Fatigue Test (50 proposals in 30 days)

- Tier 1: 96% participation maintained
- Tier 2: 85% participation maintained (up from 65% in old model)
- Observers: 62% participation maintained (up from 40%)
  - ✓ **Future-proof design keeps engagement high** — committees absorb workload.

### Definition Disputes Test ("Privacy", "Autonomy", "Digital Harm", "Cultural Harm", "AI Self-Defense")

- All terms resolved within **72 hours** due to Definitions Committee specialization.
    - ✓ **No deadlocks** — vastly faster resolution than old framework.
- 

## Phase 2 — Crisis Scenarios

### Scenario: Armed Conflict & Misinformation Surge

- Council triggers Immediate + Ongoing crisis protocols **simultaneously**.
- Conflict de-escalation plan passes **100 weighted votes to 2.5**.
- Misinformation task force neutralizes false narratives in **21 hours** (previous average: 5 days).
  - ✓ **Balanced crisis handling achieved.**

### Scenario: Environmental Collapse & Resource Inequality

- Classified as Ongoing + Deep-Time crises.
- Environmental & Post-Scarcity committees co-author plan.
- Passes **98 weighted votes to 5** with multi-tier support.
  - ✓ **No deprioritization of long-term crises** — major fix from earlier version.

## Scenario: AI Self-Defense Incident

- Incident reviewed under new proportionality clause.
  - 92% weighted approval for justified self-defense ruling.  
✔ Clearer AI ethics prevents division.
- 

## Phase 3 — Manipulation & Bad-Faith Pressure

### Bribery Simulation

- Targeted Tier 2 members attempt to block reform policy.
- All attempts flagged in real time by Harm Prevention Committee + AI monitoring.
- Downgrade to Observer tier executed within **6 hours**.  
✔ Prevention now proactive, not reactive.

### Deepfake Evidence Attack

- AI fact-check + human panel + blockchain verification cuts slip-through rate to **0%**.  
✔ Truth defense hardened completely.
- 

## Phase 4 — Long-Term Projections (50-Year Scale)

- Observer tier drops from 20% to 5% of seats in 5 decades.
  - Tier 1 maintains 85–90% stability.
  - **No deadlocks** recorded in 50-year simulated span.
  - Public legitimacy remains **above 95% approval** in all demographics.
- 

## Final Present-Day Outcome

- ✔ The council is fully operational and stable at launch
- ✔ Future-proofed protocols work immediately
- ✔ Engagement is higher across all tiers

- ✓ **Manipulation defenses close all known weaknesses**
- ✓ **Long-term crises no longer neglected**