

# Harm Blindness Framework Implementation Checklist

**Version: 2.0**

**Date:** November 19, 2025

**Created by:** Real Safety AI Foundation (Travis Gilly, executive director)

**Purpose:** This checklist guides you through implementing the Harm Blindness Framework in your organization. Work through each phase sequentially, checking off items as you complete them.

**Estimated Timeline:** 8-12 weeks from start to full rollout

---

## PHASE 1: PREPARATION

### Leadership Alignment (Week 1)

- ☐ Executive sponsor identified and committed
- ☐ Leadership team reads full Harm Blindness Framework document
- ☐ Budget allocated for implementation (time, resources, potential delays)
- ☐ Authority structure defined (who can delay/cancel based on checkpoints)
- ☐ Cultural readiness assessed (psychological safety to raise concerns)

### Team Structure (Week 1)

- ☐ Checkpoint facilitator(s) identified
- ☐ Facilitators have appropriate authority and skillset
- ☐ Project owners briefed on their responsibilities
- ☐ Stakeholder representative strategy defined
- ☐ Documentation owner assigned

### Tools & Materials (Week 1-2)

- ☐ Checkpoint templates downloaded and customized
- ☐ Stakeholder mapping tools prepared
- ☐ Risk assessment worksheets ready
- ☐ Documentation system selected (wiki, docs, PM tool)
- ☐ Templates integrated into existing tools (Jira, Notion, etc.)

### Training (Week 2)

- ☐ Key personnel complete framework training
- ☐ Facilitators trained on running checkpoints
- ☐ Q&A session held to address concerns

- ☐ Common objections documented with responses
  - ☐ Training materials prepared for broader rollout
- 

## **PHASE 2: PILOT PROGRAM**

### **Project Selection (Week 2-3)**

- ☐ 2-3 pilot projects identified
- ☐ Projects are medium-complexity (not trivial, not mission-critical)
- ☐ Project teams briefed on pilot participation
- ☐ Pilot success criteria defined
- ☐ Evaluation metrics established

### **Checkpoint Execution (Weeks 3-6)**

#### **For Each Pilot Project:**

#### **Checkpoint 1 (Ideation)**

- ☐ Checkpoint 1 completed and documented
- ☐ Death risk screening completed (Does this system pose death risk?)
- ☐ If YES: Death Gate Protocol activated (see Part 2 of framework)

#### **Checkpoint 2 (Design)**

- ☐ Checkpoint 2 completed and documented
- ☐ Death risk screening completed (Does this system pose death risk?)
- ☐ If YES: Death Gate Protocol activated (see Part 2 of framework)

#### **Checkpoint 3 (Testing)**

- ☐ Checkpoint 3 completed and documented
- ☐ Death risk screening completed (Does this system pose death risk?)
- ☐ If YES: Death Gate Protocol activated (see Part 2 of framework)

#### **Checkpoint 4 (Launch)**

- ☐ Checkpoint 4 completed and documented
- ☐ Death risk screening completed (Does this system pose death risk?)
- ☐ If YES: Death Gate Protocol activated (see Part 2 of framework)

### **Metrics Tracking**

- ☐ Time spent on each checkpoint tracked
- ☐ Issues caught documented
- ☐ Participant feedback collected

### **Pilot Evaluation (Weeks 7-8)**

- ☐ All pilot checkpoint documentation reviewed for quality
  - ☐ Time overhead calculated (should be <10% of project time)
  - ☐ Issues caught vs missed analyzed
  - ☐ ROI calculated (cost of catching issues early vs late)
  - ☐ Team satisfaction surveyed
  - ☐ Process improvements identified
  - ☐ Checkpoint questions refined based on learnings
  - ☐ Templates updated based on feedback
- 

## **PHASE 3: ROLLOUT**

### **Process Refinement (Week 8-9)**

- ☐ Checkpoint questions finalized for your context
- ☐ Documentation templates updated
- ☐ Integration with PM tools finalized
- ☐ Workflow integration documented (Agile, Waterfall, etc.)
- ☐ Success metrics defined for organization-wide tracking

### **Training Scale-Up (Week 9-10)**

- ☐ Additional facilitators trained
- ☐ All project owners trained
- ☐ All technical leads briefed
- ☐ Organization-wide announcement made
- ☐ FAQ document created and distributed

### **Mandate & Monitoring (Week 10+)**

- ☐ Framework mandated for all relevant projects
- ☐ Checkpoint completion integrated into Definition of Done
- ☐ Monitoring dashboard created
- ☐ Compliance tracking established
- ☐ Escalation path defined for non-compliance
- ☐ Success stories documented and shared

---

## **PHASE 4: CONTINUOUS IMPROVEMENT**

### **Quarterly Review Process (Ongoing)**

**Every Quarter, Review:**

#### **Metrics Dashboard**

- ☐ Are we hitting targets?
- ☐ What trends are emerging?
- ☐ Where are we struggling?

#### **Success Stories**

- ☐ What did framework catch this quarter?
- ☐ What would have happened without it?
- ☐ Document for case studies

#### **Failure Analysis**

- ☐ What harms occurred despite framework?
- ☐ What did we miss?
- ☐ How do we improve?

#### **Process Feedback**

- ☐ What's working well?
- ☐ What's frustrating teams?
- ☐ How can we streamline?

#### **Framework Updates**

- ☐ New checkpoint questions needed?
- ☐ Better templates available?
- ☐ Improved training materials?
- ☐ MIT AI Risk Repository domains reviewed as supplementary check

---

## **DEATH GATE PROTOCOL COMPLIANCE (if triggered)**

**Only required if death risk identified at any checkpoint:**

- ☐ **Stage 1: Public Warning**

- ☐ Death hazard warning implemented on all interfaces
- ☐ SEC/regulatory filing completed
- ☐ Executive personal sign-off obtained

☐ **Stage 2: Regulatory Authorization**

- ☐ Comprehensive regulatory audit requested
- ☐ Public comment period (90 days minimum)
- ☐ Legislative notification completed

☐ **Stage 3: Independent Coalition**

- ☐ Coalition of 10+ members formed
  - ☐ Supermajority approval obtained (8/10)
  - ☐ Ongoing oversight established
- 

**Notes**

- Death risk includes: Direct causation of preventable death, suicide facilitation, violence enabling capabilities, life-critical system failures
- MIT AI Risk Repository review provides comprehensive supplementary coverage across 7 domains and 24 subdomains
- Version 2.0 updates incorporate Death Gate Protocol from Part 2 of the Harm Blindness Framework