Choosing the Right Controls for AI Risks



Al Risk	Control Purpose	💡 Design-Time	🐥 Run-Time
Model Drift and Data Distribution Shift	Prevention	Use diverse and up-to-date training data	Deploy adaptive systems with online learning
	•	Apply data augmentation and simulation	Pre-schedule automatic retraining jobs
		Plan for periodic retraining	Include fallback mechanisms for high-risk scenarios
Model performance degrades over time as real-	Detection	Eulid in robustness to variation Embed drift detection hooks or components in system architecture	□ Monitor live data with drift metrics (e.g. PSI. KI_divergence)
world data diverges from training data. This can		□ Stress test with simulated distribution shifts	□ Track model performance via outcomes or feedback loops
lead to increasing errors or unfair outcomes if	Response	Define thresholds for retraining and escalation	Trigger model retraining or fine-tuning
ongoing monitoring.	·	Design retraining pipelines into architecture	Rollback to stable models if needed
			Escalate to human review or rules-based system in critical cases
Hallucinations in Generative Models	Yrevention	Fine-tune on high-quality, verified, domain-specific data	Apply confidence thresholds to suppress or disclaim low-certainty outputs
		Use Retrieval-Augmented Generation (RAG) to ground responses in trusted sources	Use known prompt handling rules to avoid hallucination traps
		Apply Reinforcement learning from human Feedback (RLHF) Design abstention mechanisms (e.g. "I'm not sure" responses)	□ Use human-in-the-loop review in high-risk contexts
Generative models can produce plausible-	Detection	Train with human feedback to improve truthfulness metrics	Run real-time fact-checking on generated content against trusted sources
sounding but false or fabricated outputs. These	~	□ Build in self-consistency or ensemble checking mechanisms	Use secondary models or filters to evaluate output accuracy
hallucinations can mislead users, spread			Allow user reporting of suspected hallucinations
Controls are needed to promote truthfulness	Response	Incorporate fallback logic for uncertain outputs (e.g. human/structured answer)	Flag or block hallucinated content
and catch inaccuracies during generation.			Escalate to human review if fact-checking fails
			Log and learn from natiocination cases to improve future accuracy
Bias and Fairness Issues	Prevention	Curate diverse and representative training datasets	Enforce fairness rules in real-time (e.g., content balancing in recommendations)
Δ	•	Apply fairness-aware algorithms (e.g., re-weighting, fairness-constrained optimization)	Use human-in-the-loop to review high-impact decisions
		 Define and validate against fairness metrics (e.g., demographic parity, equalised odds) Conductions do logical straight fairness metrics (e.g., demographic parity, equalised odds) 	Apply dynamic fairness constraints based on live inputs and user demographics
Al systems may exhibit or reinforce bias leading		Conduct pre-deployment audits and fairness stress tests Include diverse perspectives in design reviews	
to unfair or discriminatory outcomes. This can	Detection	Simulate outcomes for different groups during validation	Continuously monitor outcomes by group (e.g., acceptance rates by race/gender)
result from biased training data, model design,		□ Audit model performance across subpopulations	 Perform regular fairness audits on system outputs
or unintended feedback loops. Addressing this risk is essential to ensure ethical, legal, and		Identify proxy variables that may encode bias	Track feedback loops or data shifts that may introduce bias over time
reputational integrity.	Response	Establish remediation plans for failed fairness tests	Provide explanations and recourse for affected users
		Iterate model or feature set to reduce disparate impact	Retrain or adjust models when bias is detected in operation
			□ Trigger numan review or escalation protocols for hagged cases
Adversarial Inputs and	🔶 Prevention	Use adversarial training with attack examples	Apply rate limiting and throttling to prevent attack probing
Robustness Vulnerabilities	•	□ Preprocess inputs to remove adversarial noise (e.g., input normalization)	Sanitise or validate user inputs before model access
N		Perform red-team testing and security reviews	Limit access to sensitive model capabilities (e.g., restricted prompts)
Malicious actors may exploit Al systems using		Design ensemble or redundant model architectures	
crafted inputs or attacks (e.g., adversarial	Detection	Simulate adversarial scenarios during evaluation	Monitor input patterns for anomalies or adversarial characteristics
examples, prompt injections, data poisoning).		 Test models with known adversarial patterns 	Use confidence-based or activation-based anomaly detectors
These attacks can bypass or mislead models, especially in security-critical applications		·	□ Log inputs and outputs for forensic analysis and threat pattern recognition
Defences must anticipate, detect, and respond	Response	Plan recovery paths for adversarial failure modes	Trigger alerts or shutdown mechanisms if attack is suspected
to such threats.		Embed fallback models or decision logic for high-risk inputs	Isolate or block malicious inputs in real time
			Patch models or filters in response to discovered vulnerabilities Engage AI security teams for live response and threat mitigation
			a Engage / recently count of the response and threat integration
Loss of Personal or	Prevention	Apply data minimisation: exclude unnecessary PII	Filter outputs in real-time for PII or sensitive patterns
Confidential Information		Audit training datasets for sensitive content	Limit access to models and outputs via authentication and roles
		Ose differential privacy or regularisation to reduce memorisation	Use input sanitisation to block confidential user submissions Ware users not to share sensitive information in prompts
AI systems may inadvertently expose sensitive,		 Anonymise of mask data before training Define strict access roles and separation of duties in system design 	• Warn users not to share sensitive information in prompts
personal, or internal data through	Detection	Test for memorisation using known PII probes	Use automated output scanning (NER, PII detectors)
memorization, output generation, or insecure	_	Run red-team attacks to surface potential leakage	Log and monitor for suspicious output or access behavior
concerns – especially under data protection		Monitor training for overfitting to rare or personal data	Conduct audits for data exposure or anomalies
laws.	Response	Establish privacy risk thresholds for retraining or model blocking	Takedown or suppress harmful outputs
		Document potential leakage risks in system governance reviews	 Notity affected users or regulators if required Trigger incident response protocols for potential breaches
			 Update training and filtering strategies based on incidents
Harmful Content	Prevention	Implement strict data curation practices to evolude harmful content during training	Employ real-time content filtering systems to detect and block barmful outputs
Generation or Exposure		Apply content moderation policies aligned with legal and ethical standards	 Implement user input sanitisation to prevent the generation of harmful outputs
		 Design models with built-in safety constraints to limit the generation of harmful content 	Establish user reporting mechanisms for harmful content
		Use Reinforcement Learning from Human Feedback (RLHF) to limit undesirable outputs	Apply rate limiting and monitoring to detect and prevent abuse patterns
Al systems, especially generative models, may	Oetection	Conduct adversarial testing to identify potential for harmful content generation	Monitor content outputs continuously for signs of harmful material
offensive, abusive, or otherwise harmful. This		Develop classifiers to detect harmful content within model outputs Establish banchmarks for accordable content and test are delegative these stands are	Utilize automated tools to detect and flag harmful content in real-time
includes hate speech, explicit material,	Response	Establish benchmarks for acceptable content and test models against these standards Create protocols for content takedown and user notification	Analyze user recovack and reports to identify patterns or narmful Content generation Implement immediate content removal or correction mechanisms
misinformation, or content that incites violence.	/	 Develop incident response plans for addressing harmful content-related crises 	 Suspend or modify Al system functionalities when harmful content is detected
Such outputs can lead to user harm, legal repercussions, and reputational damage.			Engage human moderators to review and address flagged content promptly
			Update models and filters based on incidents to prevent future occurrences
Feedback Loops and Behaviour Amplification	Prevention	□ Introduce diversity-promoting mechanisms in algorithms (e.g., explore-exploit strategies)	Periodically inject randomised or diverse content to break loops
	*	Apply multi-objective optimization (e.g., engagement + diversity + well-being)	Allow users to actively request content variation ("show me something new")
		Conduct simulation studies to test long-term behavior patterns	Reset or perturb model states on a scheduled basis
		Enable user control features (e.g., reset or broaden profile options)	
Al systems can unintentionally reinforce behaviors or biases by acting on and shaping	Detection	Set policy constraints to prevent repetitive or overly narrow recommendations Run simulations during development to observe feedback effects	Monitor matrice like content diversity and again and concentration or user palarization over time
		en annuations during development to observe recuback Effects	=

their own input data — creating self-reinforcing cycles. This can lead to echo chambers, polarisation, or instability in dynamic environments like social media or financial markets.	Response	 Analyse how model behavior evolves across interaction cycles Tune algorithms based on simulation findings Embed decay factors or reset mechanisms into the system logic 	 Detect signs of degenerative cycles or user behavior anomalies Audit system impact on group behaviors and preferences Adjust model parameters in real time to counter amplification trends Escalate to human review when feedback effects exceed thresholds
			Periodically retrain on external or unbiased data to rebalance the system
Overreliance on Automation / Erosion of Human Oversight	Prevention	 Design transparent, explainable AI outputs Include cognitive forcing functions (e.g., rationale for accepting AI recommendations) Provide alternative perspectives (e.g., second opinion models) Train users during rollout on AI limitations and failure cases 	 Include mandatory review steps for high-risk recommendations Provide in-system prompts or nudges to encourage critical evaluation Limit duration of fully automated operation before requiring human check-in
As AI systems become more capable, users may place too much trust in them, leading to passivity or failure to challenge incorrect outputs. This "automation bias" is especially dangerous in high-stakes domains like healthcare, finance, or aviation. Effective controls preserve human judgment and ensure the AI is treated as a decision aid, not a final authority.		 Build systems that defer to humans when uncertain Establish rules requiring human sign-off for critical decisions 	
	 Detection 	 Evaluate interface design for undue trust tendencies Test user workflows for blind acceptance of AI suggestions 	 Monitor user behavior for signs of overreliance (e.g., always accepting Al outputs without edits) Track model confidence and flag low-certainty cases for human review Conduct regular oversight audits and operator feedback reviews
	🗪 Response	 Adjust system design based on user testing outcomes Reinforce human-in-the-loop workflows in high-risk contexts 	 Trigger alerts or warnings when overreliance is detected Escalate questionable outputs for human validation Share examples of caught AI errors to reinforce vigilance culture Retrain staff or adjust incentive structures to reward thorough review

 ${\it Subscribe for more free resources @ Doing Al Governance www.ethos-ai.org}$