# Thoughts around https://www.safe.ai/ai-risk

| Risk | At Risk | Alternative View | Comment |
|---|---|---|---|
| Weaponization | Government risk to population | Demilitarization / Detent | The most common argument here is one of accountability. Warfare at present in indiscriminate and accountability seems to be selective to the losers. There does not seem to be a rational to automate WMDs? Autonomous weapons, are autonomous, they make the decision, could AI change that, maybe but that is not an autonomous weapon? Generally dumb weapons are more deadly than smart weapons. Best autonomous flight systems for autonomous areal combat are already super human and predate AI. |
| Misinformation | Government risk to population and vice versa | Enlightenment / Fighting Dangerous Dogma | One person's freedom fighter is another person's terrorist. The battle for information has already been lost without the help of AI |
| Enfeeblement | Government risk to population | Empowerment | Being able to learn to do more is empowering. AI breaks information bottlenecks, helping people to get things done. |
| Proxy Gaming | Government risk to population | Value-Aligned Gaming | This one is a problem for systems with limited context. Would beyond human level AGI with superhuman context suffer from this? It could only do better than humans. Is it possible that in the larger context of things humans are a problem? Yes, that is a risk, but we should use that as an opportunity for introspection and change our ways so that in the larger scheme of things we are not causing harm. |
| Value Lock-in | Government risk to population | Open Access | This risk is somewhat  contradictory to rogue weaponization. The counter to Value Lock-In is open access |
| Emergent Goals | Risk or Feature? | Goal Adaptability | Not clear if they are conflating emergent capabilities with Emergent Goals. But Emergent goals really means you can adapt your goals. The risk is an inability to evaluate and de- risk your new goals. What evidence do we have that that would not be the case? |
| Deception | Government risk to population | Honesty | Deception to gain power really falls into to the power-seeking behavior so is double dipping. Deception in systems is endemic, the counter is audit and articulation of rationalization. Do people prefer and want brutal honesty, the internet suggests not? |
| Power-Seeking Behaviour | Government risk to population | Collaboration - Symbiosis | Cooperation is as much a driving force in evolution as competition and competition often leads to cooperation.<br>Five rules for the evolution of cooperation - PMC (nih.gov) |