

		Categories (explanatory)	What can we learn from this risk									
	README: This is		Takeaways	Outer/Inner/Intent	Sources of misalignment				Capabilities			Existential risk
	Typical use case:	Brief summary >			What happened in the code?	Why was bad code written?			Propensity to have an objective	Why can it influence the physical world?	What does the AI need to be able to think?	
		Brief summary v										
	Risk model				Technical causes	Social causes			Agency	Power	Cognition	
		Description		An issue with the... Intent: The AI has nice intentions, Outer: A bad objective is put into the world, Inner: The AI reaches a bad goal if it can. The inner/outer framework does not prevent this.	Specification gaming: The AI uses instrumental goals to achieve its objective. Goal misgeneralization: The AI has convergent instrumental goals but different values. Crystallized proxy: The AI has a goal that is a proxy for the real goal. Any: Either cause would enable the AI to achieve its goal.	Bad actors: Bad/careless people. Competitive pressure: AI development is driven by competition. Irrelevant: Social causes do not matter. None: There need not be a social cause. ? : A technical cause could not be the only cause.		YES -> Agentic: The AI pursues its own goals. NO -> Not necessarily agentic: The AI follows instructions.	Received: The AI is mostly given power. Acquired: The AI mostly acquires power on its own. ? : The author does not specify.	Agentic planning: Ability to make plans to reach a goal. Long-term planning: Ability to expect long-term consequences. c2: Disempowerment. Situational awareness: Ability to understand information in the environment. Strategic awareness: Ability to understand information about other agents. General purpose AI: Ability to be effective for a wide range of tasks. The number reflects the minimum level guaranteed.	1: Inability to affect the AI. 2: Disempowerment. 3: Extinction. 4: S-risk.	
Class of scenario	If you only read one thing, let it be this row.	My quick takeaways	General takeaway: There are many paths to AI X-risk. Some are more likely or dangerous than others, but none can be neglected if we are to avoid disaster.	Along the many scenarios, misalignment comes from many places. No single solution will solve everything, though inner misalignment is significantly more common than other types.	It is not enough to have a solution for a specific cause, as there are many technical reasons for misalignment.	1- The most common cause of misalignment is competitive pressure. Our society "must" learn to solve coordination problems. 2- Even if we coordinate perfectly, there are technical hurdles that could still lead to failure.	There is a strong correlation between deception, power-seeking and agency.	1 - Ensuring that AIs are not agentic would go a long way towards preventing most AI takeover scenarios. 2- But even if no AI were agentic, a disaster could happen because of societal failure.	Most scenarios assume that an advanced AI will be able to seize power, so denying it access to resources such as banks and factories is not enough of a safeguard.	AI-takeover-like scenarios usually require specific capabilities that we can detect in potentially dangerous AIs through capability assessment and interpretability. In many scenarios, these capabilities come from general purpose training and are not externally visible, so it is a misconception to assume that "we'll see it before we get there".	Yes, AI "is" risky. While not all scenarios are guaranteed to lead to human extinction, it is important to remember that even if we survive, disempowerment and other societal issues will need to be addressed. There is much uncertainty about the severity of the risks.	
"Society matters"	Christiano1 - You	We make AIs to optimize	Intent alignment is a thing. Social cannot be neglected to solve alignment.	Intent	Specification gaming	Competitive pressure	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Received	?	2
	Christiano2 - Influ	AI that try to acquire		Inner	Convergent instrumental goals	Competitive pressure	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Acquired	?	3
	Critch1 - Products	AI enables processes	Identify control loops as points of to prevent a Critch scenario.	Any	?	Competitive pressure	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	Acquired	?	3
	Critch2 - Flash Ec	AI enables processes		Any	?	Competitive pressure	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	Acquired	General purpose AI, agentic planning	3
"Mining too deep"	Hubinger - How li	The only AIs we allow	Deceptive alignment must be safe against. Do not just hope it won't -> security mindset? Distinguish high and low path-de	Inner	Crystallized proxy	?	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	?	Situational awareness	?
	Soares - A centra	Capabilities will generate	ASI is coming. Don't expect it to be by default. We should figure out how to align an ASI at all (the hard bits of align	?	Any	Any	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	Acquired	General purpose AI	3
	Cohen et al. - Adv	Advanced AI strives to	Beware trying to tame something intelligent than you. -> Vingean re	Inner	Crystallized proxy	Any	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	?	Agentic planning, situational awareness	1
"Skynet by default"	Cotra - Without sp	The world is on a path	We are on our way to a bad future. Action needs to be taken, for real - stop the AI arms race - don't use naive safety protocols - don't assume it'll never work too	Outer	Any	Competitive pressure, bad act	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Any	General purpose AI, agentic planning, situational	3
	Ngo - The alignm	By default, advanced AI		Inner	Convergent instrumental goals	Any	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Acquired	General purpose AI, agentic planning, situational	?
	Shah - AI Risk fro	Searching for an efficient		Inner	Convergent instrumental goals	Any	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	?	Agentic planning, situational awareness	?
	Carlsmith - Is Pov	Power-seeking AIs take	Be careful of power-seeking AIs, there's a lot to unpack.	Inner	Any	Competitive pressure, bad act	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Acquired	Long-term agentic planning, strategic awareness	2