

	<p><b>All comments will be made public as-is, with no edits or redactions. Please be careful to not include confidential business or personal information, otherwise sensitive or protected information, or any information you do not wish to be posted.</b></p>					
			<p>Comment Template for Responses to NIST Artificial Intelligence Risk Management Framework Request for Information (RFI)</p>			<p><b>Submit comments by August 19, 2021:</b></p>

<b>General RFI Topics</b> (Use as many lines as you like)	<b>Response #</b>	<b>Responding organization</b>	<b>Responder's name</b>	<b>Paper Section (if applicable)</b>	<b>Response/Comment (Include rationale)</b>	<b>Suggested change</b>
			Aryeh Englander		NIST needs to be proactively reaching out to relevant research communities and organizations, especially the AI safety research community, and the assured autonomy / AI safety engineering community. In particular, NIST should be proactively talking with the following organizations: Assuring Autonomy International Programme ( <a href="https://www.york.ac.uk/assuring-autonomy/">https://www.york.ac.uk/assuring-autonomy/</a> ); Center for Human-Compatible Artificial Intelligence ( <a href="https://humancompatible.ai/">https://humancompatible.ai/</a> ); Center for Security and Emerging Technology (CSET, <a href="https://cset.georgetown.edu/">https://cset.georgetown.edu/</a> ); Anthropic ( <a href="https://www.anthropic.com/">https://www.anthropic.com/</a> ); Stanford University Human-Centered Artificial Intelligence (HAI, <a href="https://hai.stanford.edu/">https://hai.stanford.edu/</a> ); Consortium on the Landscape of AI Safety (CLAIS, <a href="https://www.clais.org/">https://www.clais.org/</a> )	
			Aryeh Englander		The following references seem very relevant, and NIST should be incorporating the issues and suggestions discussed in these references: <a href="https://arxiv.org/abs/1905.04223v1">https://arxiv.org/abs/1905.04223v1</a> , <a href="https://www.york.ac.uk/assuring-autonomy/guidance/amlas/">https://www.york.ac.uk/assuring-autonomy/guidance/amlas/</a> , <a href="https://arxiv.org/abs/2004.07213v2">https://arxiv.org/abs/2004.07213v2</a> , <a href="https://www.york.ac.uk/assuring-autonomy/guidance/body-of-knowledge/">https://www.york.ac.uk/assuring-autonomy/guidance/body-of-knowledge/</a> , <a href="https://arxiv.org/abs/2108.07258v2">https://arxiv.org/abs/2108.07258v2</a> , <a href="https://cset.georgetown.edu/publication/ai-accidents-an-emerging-threat/">https://cset.georgetown.edu/publication/ai-accidents-an-emerging-threat/</a>	
			I-Jeng Wang		There is a need to imbue risk-sensitive behavior into advanced AI agent to enable shaping of risk-taking strategies consistent with human values.	

<b>Responses to Specific Request for information</b> (pages 11,12, 13 and 14 of the RFI)						
1. The greatest challenges in improving how AI actors manage AI-related risks – where “manage” means identify, assess, prioritize, respond to, or communicate those risks;			Aryeh Englander		Difficult or impossible to test all possible situations including edge cases for very complex AI systems deployed in very complex environments.	
			Aryeh Englander		Many complex AI systems, especially deep learning networks, are essentially black boxes which are extremely difficult or impossible to understand with current techniques.	
			Aryeh Englander		Advanced machine learning systems are liable to discover very novel solutions that satisfy the objective we gave it but which may not satisfy what we actually want. For very advanced systems it becomes extremely difficult or impossible to precisely specify everything that we do or do not want the system to do, which could lead to the AI finding novel solutions that we very much do not want, in ways that we may not know about until it is too late to prevent.	
			Aryeh Englander		There is currently very little serious government-level discussion of globally catastrophic or even existential risks from very advanced AI systems, despite warnings from many experts that very advanced AI may pose such threats within the next few decades. Decision makers often dismiss such concerns as "science fiction" without actually looking at the relevant arguments and evidence.	

			Aryeh Englander		There is a very difficult challenge of setting up auditing, oversight, and governance mechanisms so that actors and organizations actually follow through with the principles they say they agree to. Very often organizations create lists of great principles that they will adhere to, and then those principles end up being mostly a PR piece and they only get very poorly instituted in practice if at all.	
			Aryeh Englander		Technology races between companies or nations is perhaps the greatest factor in AI risks. If implementing safety or ethical concerns turns out to be difficult or costly or if it negatively impacts performance, then a "race to the bottom" becomes highly likely, where competing companies or nations are greatly incentivized to cut corners in terms of safety or ethics. National and international risk mitigation frameworks, perhaps including treaties between nations, may be critical for solving this issue.	
			James P. Howard		The greatest challenge is in quantifying risk. Some risks are easily quantified because we know the associated risk of error. However, some risks are unquantifiable due to rarity, lack of data, or lack of knowledge the risk even exists. Catch-all risk management can attempt to capture this, but it is hit or miss at best.	
			Katie Zaback		<p>Lack of quantification of uncertainty/risk (i.e. current AI systems often do not have quantified representation of risk/uncertainty baked into the algorithm - it's either non-existent/secondary/"soft")</p> <p>Real-world risk is difficult to quantify; even more difficult - quantifying how to determine what is an acceptable "level" of risk [this is usually domain/application specific]</p> <p>Risks can be introduced not just in results/predictions but in blind-spots, inherent bias or "invisible" risk that might be baked into the underpinnings of the algorithm or that data that drives the algorithm</p> <p>Related: Neural networks (and similar methods) will always carry risk of error (i.e. failure/incorrect predictions/etc) will ALWAYS be a part of the system.</p>	
			I-Jeng Wang		It is extremely difficult to predict risks associated with dynamic context-dependent adaptation of envisioned online or lifelong learning AI. This is especially concerning due to learned optimization. See <a href="https://arxiv.org/abs/1906.01820">https://arxiv.org/abs/1906.01820</a> .	

<p>2. How organizations currently define and manage characteristics of AI trustworthiness and whether there are important characteristics which should be considered in the Framework besides: accuracy, explainability and interpretability, reliability, privacy, robustness, safety, security (resilience), and mitigation of harmful bias, or harmful outcomes from misuse of the AI;</p>			I-Jeng Wang		<p>Uncertainty estimates on output of advanced ML techniques such as DNN shall be a key elements of any AI models deployed to safety critical domains. A comprehensive and effective uncertainty modeling framework/methodology is lacking and remains an open research problem.</p>	

<p>3. How organizations currently define and manage principles of AI trustworthiness and whether there are important principles which should be considered in the Framework besides: transparency, fairness, and accountability;</p>			<p>Katie Zaback</p>		<p>Trustworthiness should be more granular than high-level markers of transparency and fairness. Details assessments related to data would be a good place to start (i.e. describe the data that was and was not present when training; this should be more than data on the <i>*amount*</i> of data, but also the <i>*type*</i> and <i>*quality*</i>).</p> <p>Trustworthiness also relates to human perceptions. For example, a pilot that is in a plane being controlled in part by an AI system might not trust that system with her life in practice, but would in a virtual setting. Thinking about these issues of trust when building (and implementing) these AI systems should be considered. (Good example: DARPA's Alpha Dogfight)</p>	
<p>4. The extent to which AI risks are incorporated into different organizations' overarching enterprise risk management – including, but not limited to, the management of risks related to cybersecurity, privacy, and safety;</p>			<p>Aryeh Englander</p>		<p>Government organizations and other large-scale organizations need to be actively incorporating longer-term considerations of risks from very advanced AI that experts anticipate may be coming in the next few decades. There is considerable research that can be done now to mitigate those risks, yet very few organizations are thinking about them.</p>	

			Katie Zaback		Again, this is application/domain specific. Some AI systems might be highly correlated with safety, but not with privacy. Some might represent significant cybersecurity risks, but not concerns of physical safety. This presents difficulty when building a robust framework for risk.	
5. Standards, frameworks, models, methodologies, tools, guidelines and best practices, and principles to identify, assess, prioritize, mitigate, or communicate AI risk and whether any currently meet the minimum attributes described above;			Aryeh Englander		NIST may also want to consider the following risk framework for decisions related to longer-term AI risks: <a href="https://www.alignmentforum.org/posts/qnA6paRwMky3Q6kktk/modelling-transformative-ai-risks-ntair-project-introduction">https://www.alignmentforum.org/posts/qnA6paRwMky3Q6kktk/modelling-transformative-ai-risks-ntair-project-introduction</a> (the author of this comment is a POC for this project)	
6. How current regulatory or regulatory reporting requirements (e.g., local, state, national, international) relate to the use of AI standards, frameworks,						

models, methodologies, tools, guidelines and best practices, and principles;						
7. AI risk management standards, frameworks, models, methodologies, tools, guidelines and best practices, principles, and practices which NIST should consider to ensure that the AI RMF aligns with and supports other efforts;			James P. Howard		NIST should look at the economic research on decision theory, especially decision-making under uncertainty. This is used in the economic field to make decisions without complete information. This is, conceptually, isomorphic to the question posed by NIST here. Risk management standards around uncertainty can be used to support how an organization responds to a risk (see, for instance, how an investment bank makes investment decisions).	

<p>8. How organizations take into account benefits and issues related to inclusiveness in AI design, development, use and evaluation – and how AI design and development may be carried out in a way that reduces or manages the risk of potential negative impact on individuals, groups, and society.</p>			<p>Katie Zaback</p>		<p>Reframe "inclusiveness" as "anti-discriminatory measures." Inclusiveness downplays the harm that will be done by AI systems that are improperly trained. Make anti-discriminatory practices when designing/training/deploying AI systems a requirement of a good system - not something secondary to the design of the system. To do otherwise doesn't clearly communicate that a discriminatory system is not only harmful, but does not meet system requirements. For example, a facial recognition system that performs poorly on black faces is not a facial recognition system - it is a *white* facial recognition system (i.e. does NOT meet system requirements).</p>	
<p>9. The appropriateness of the attributes NIST has developed for the AI Risk Management Framework. (See above, "AI RMF Development and Attributes");</p>						

<p>10. Effective ways to structure the Framework to achieve the desired goals, including, but not limited to, integrating AI risk management processes with organizational processes for developing products and services for better outcomes in terms of trustworthiness and management of AI risks. Respondents are asked to identify any current models which would be effective. These could include – but are not limited to – the NIST Cybersecurity Framework or Privacy</p>			<p>James P. Howard</p>		<p>NIST could propose a maturity model for risk management under AI decision-making. For instance, at the low end of a maturity model, there is no risk management. At the high-end can be a complete response framework built into an organizational decision-making. By defining as a maturity model, an organization can adapt the requirements to their specific needs, even on a project-by-project basis.</p> <p>In addition, NIST can again look toward the financial community to structure these responses and look toward the Basel Accords for banking risk management and the Solvency II directive issued by the EU for insurance risk management.</p>	

Framework, which focus on outcomes, functions, categories and subcategories and also offer options for developing profiles reflecting current and desired approaches as well as tiers to describe degree of framework implementation; and						

<p>11. How the Framework could be developed to advance the recruitment, hiring, development, and retention of a knowledgeable and skilled workforce necessary to perform AI-related functions within organizations.</p>						
<p>12. The extent to which the Framework should include governance issues, including but not limited to make up of design and development teams, monitoring and evaluation, and grievance and redress.</p>						