

Article

# AI Governance and the Policymaking Process: Key Considerations for Reducing AI Risk

Brandon Perry <sup>1,\*</sup>, Risto Uuk <sup>2</sup>

<sup>1</sup> Independent researcher; brandonperryofficial@gmail.com

<sup>2</sup> Effective Altruism Estonia; ristouuk@gmail.com

\* Correspondence: brandonperryofficial@gmail.com

Version April 4, 2019 submitted to Big Data Cogn. Comput.

**Abstract:** In recent years, a growing number of researchers in a variety of fields have begun working on the question of how to mitigate the catastrophic risks that could result from transformative artificial intelligence, including what policies states should adopt in response. However, in this paper we identify a preceding, meta-level problem of how the space of possible policies is affected by the policymaking process itself, which includes: (1) problem identification/agenda setting, (2) policy formulation, (3) policy adoption, (4) implementation, and (5) evaluation. The features of the policymaking process not only potentially constitute a large obstacle or opportunity in the way for solving AI risk, but it also creates a new set of key considerations for the field at large. In this paper, we argue that a new subfield of AI governance should be explored that examines the policy-making cycle and its implications for AI governance, and then end the paper with key considerations to guide future areas of research.

**Keywords:** policymaking process; AI risk; typologies of AI policy; AI governance

## 1. Introduction

Artificial intelligence, especially artificial general intelligence (AGI), has the ability to dramatically impact the future of humanity.[1] Notable researchers, such as Bostrom (2014), have expressed concern that advanced forms of artificial intelligence, if not aligned to humans values and wellbeing, could be potentially disastrous and pose an existential threat to our civilization.[2] The two main branches of research on risk from advanced AI are AI safety, which seeks to ensure that advanced AI is engineered in such a way that it won't pose a threat; and AI governance, which focuses on political and social dynamics (AI macrostrategy) and forecasting timelines for AI development.[3] Issues that AI governance looks at includes arms race dynamics, social and economic inequality, public perceptions, issues in surveillance, and more.

There has been a modest amount of work on developing policy solutions to AI risk, with a recent literature review by Baum (2017)[4] and Everitt (2016)[5] covering most of it. Some authors have focused on the development of AGI, with proposed solutions ranging from Joy (2000)[6] who calls for a complete moratorium on AGI research, to Hibbard (2002)[7] and Hughes (2007)[8] who advocate for regulatory regimes to prevent the emergence of harmful AGI, to McGinnis (2010) who advocates for the US to steeply accelerate friendly AGI research.[9] Everitt et al. (2017)[10] suggests that there should be an increase in AI safety funding. Scherer (2016)[11] however, at least in the context of narrow AI, argues that tort law and the existing legal structures, along with the concentration of AI R&D in large visible corporations like Google, will provide some incentives for the safe development of AI. Guihot et al. (2017)[12] also notes that attempts to future-proof laws tend to fail, and pre-emptive bans and regulation tends to hurt the long-term health of the field, instead arguing for a soft-law approach. Other authors have focused on the community of researchers, with Baum (2017)[13] promoting a social

35 psychology approach to promote community self-regulation and activism, and Yampolskiy and Fox  
36 (2013)[14] advocating for review boards at universities and other research organizations.

37 Some authors have advocated for an international approach to resolving AI risk. Erdelyi and  
38 Goldsmith (2018)[15] advocates for an international soft-law regime that would serve as a "international  
39 forum for discussion and engage in international standard setting activities." Erdelyi and Goldsmith's  
40 proposal however is not targeted towards AGI risk, although they could scale up to AGI. Wilson  
41 (2013)[16] and Bostrom (2014)[17] on the other hand call for some form of international agreement or  
42 control on AGI R&D, with the former advocating specifically for a treaty.

43 These approaches are necessary given some of the risks, including states pursuing AGI for  
44 unprecedented military and economic strength with destabilizing effects (Shulman 2009)[18], and the  
45 concentration of wealth and political influence in large corporations (Goertzel 2017).[19] Questions  
46 regarding whether or not AGI R&D should be open sourced or not have been explored by Goertzel  
47 (2017)[20] and Bostrom (2017).[21] Shulman (2009)[22], and Dewey (2015)[23] take a completely  
48 different approach, and advocate for a global surveillance regime to monitor for rogue AGI projects,  
49 with Goertzel (2012)[24] suggesting that a limited form of AGI could do this.

50 As far as current and future research goes, the Future of Humanity Institute has developed an  
51 extensive research agenda[25] for AI governance, with three main research areas: technical landscape,  
52 which seeks to understand what artificial intelligence can do and its limits; AI politics, which looks at  
53 the political dynamics between firms, governments, publics, etc; and ideal governance, which looks  
54 at possible ways and arrangements for stakeholders to cooperate. This research agenda highlights  
55 key issues such as security challenges, international political dynamics and distribution of wealth,  
56 and arms race dynamics. Other researchers have published reports dealing with issues such as dual  
57 use, similarity and possible interactions with the cybersecurity community,[26] the role and limits of  
58 principles for AI ethics,[27] justice and equity,[28] and AGI R&D community norms.[29]

59 Thus far, much of the literature on AI risk has discussed policy issues, but few have talked about  
60 how policies are made or how the dynamics of the policymaking process affect their work. Calo  
61 (2017)[30] touches upon the problem, noting that there's a lack of institutional expertise, policy tools,  
62 and flawed mental models of what AI is, which plague governments' abilities to regulate AI. Scherer  
63 (2016)[31] cites certain aspects of the technology itself, such as its ability to be created without special  
64 equipment, as a hindrance to the ability to regulate it. Everitt et al. (2017)[32] also briefly discusses  
65 policy and political dynamics in the context of AGI researchers, suggesting that AGI researchers should  
66 work with other organizations to mitigate the negative dynamics of framing AGI development as an  
67 arms race.[33] Finally, the Future of Humanity Institute's research agenda for AI governance touches  
68 on policymaking in a few ways, noting that public opinion can have major impacts on technology  
69 policy, governance schemes can be subject to mission drift, and asking how to facilitate the transition  
70 from the present state of affairs to our ideal vision for the future.

71 In this paper we will continue along the lines of facilitating the transition from the present state to  
72 our ideal vision by exploring the missing discussion on the role of policymaking in AI Governance.  
73 Research thus far has largely focused on what problems are out there and what we should do to fix  
74 them. However, throughout this paper we will not only argue that proposal implementation that takes  
75 into account the features of the 'policymaking cycle' may be vital to success in reducing AI risk, but  
76 that this model actually has massive implications for the research field as a whole. Proposals will  
77 be much more effective if they are informed by an understanding of the political and administrative  
78 considerations of consensus-building and implementation, and could make the difference between  
79 making an impact or none at all.

80 The goal of this paper is to attempt to create a launching point for discussions on the key  
81 considerations of the policymaking process for AI governance and the political considerations  
82 underpinning policy solutions for AI risk. The policymaking process includes: (1) problem  
83 identification/agenda setting, (2) policy formulation, (3) policy adoption, (4) implementation, and  
84 (5) evaluation. Each step of the policymaking process will have different aspects that are critical for

85 the creation of public policies that are able to effectively reduce AI risk. Each section will do a brief  
86 overview of the literature, assess its implications for the greater AI governance field, and identify  
87 different points where further research is needed.

88 The first section will map out and define terms in the field of AI governance, to give readers a  
89 better understanding of how our paper contributes to the way AI Governance is approached. We  
90 also create a typology for AI risk policies, to provide an understanding as to how AI governance  
91 has implications in a diverse range of policy communities, and how that interplays with strategic  
92 considerations. The next section goes through each step of the policymaking cycle, doing a basic  
93 overview of some of the literature and discussing its implications for AI governance. We would also  
94 like to note that the literature we cover in each field is not extensive and further research may be  
95 necessary. The last sections cover some of the key implications and limitations.

## 96 2. Terms and Definitions

97 On a broad level, the question of mitigating AI risk, or risks that stem from the development and  
98 use of artificial intelligence (such as global catastrophic risks from misaligned AI or military instability  
99 from adopting new types of weapons), is broken down into AI technical safety and AI governance. AI  
100 technical safety focuses on solving the computer science problems around issues like misalignment  
101 and the control problem for AGI.[34] AI governance on the other hand, studies how humanity can  
102 best navigate the transition to advanced AI systems.[35] This would include the political, military,  
103 economic, governance, and ethical considerations and aspects of the problem that advanced AI has on  
104 society.

105 AI governance can be further broken down into other components, namely the technical landscape  
106 (how technical developments depends on inputs and constraints, and affects rates or domains of  
107 capability improvement), ideal governance (what would we do ideally if we could cooperate), and  
108 AI politics (how AI will affect domestic politics, political economy, international relations, etc).[36]  
109 From these research areas, we can define the problems and solutions necessary to discuss AI policy.  
110 Throughout the paper however, we will refer to this as AI risk policy to differentiate policies intended  
111 to reduce catastrophic risk to society versus policies that apply to AI in any other circumstances.

112 Policies however, must be implemented into the legal statutes of government in order to work.  
113 Flynn (2017)[37] in the blog post that defines 'AI strategy'[38] also defines 'AI policy implementation',  
114 which is carrying out the activities necessary to safely navigate the transition to advanced AI systems.  
115 This definition implies it is action-oriented work done in government, policy, lobbying, funding,  
116 etc. As mentioned in the endnotes of Flynn (2017) however, there's an implicit gap between AI  
117 strategy (governance) research and policy implementation, with no AI policy research that identifies  
118 mechanisms for actualizing change.

119 However, there is another gap that this paper intends to address, which is that the processes that  
120 creates and implements policies (the policymaking process) often either distorts the original policy, fall  
121 short of or even work counter to the intended outcome, or render certain policy options unactionable.  
122 Similarly, The AI Governance: A Research Agenda report has neither this consideration nor a definition  
123 of policy implementation. We intend for our definition of AI policymaking strategy to fill this gap,  
124 which we define as:

125 AI Policymaking Strategy: A research field that analyzes the policymaking process and draws  
126 implications for policy design, advocacy, organizational strategy, and AI governance as a whole.

127 This goes further than the concern listed in the endnotes and also develops an upstream approach  
128 to AI governance, where work in implementation in turn feeds back, and can provide new insights to  
129 AI governance research.

130 AI policymaking strategy would fit under the definition of AI governance and would be its own  
131 sub-field in the same way technical landscape is, and would help clarify questions and considerations  
132 in the other subfields. AI politics and ideal governance seems to ask questions about what risks  
133 humanity faces and what we ought to do about them, approaching the world as if from above

134 and making corrections. Whereas policymaking strategy asks questions about how and what can  
135 be done, given both present and future circumstances, and the methods to do so at hand. They  
136 approach the world as agents who individually influence the trajectory of the world. These two groups,  
137 when they work together, should ideally converge on a policy programme that both works and is  
138 pragmatic—constituting of policies that both aim at the correct goals, and can actually get there.

139 An example of this would be the proposed solution by Goertzel (2012)[39] of creating a  
140 surveillance artificial narrow intelligence that monitors the world to prevent the development of  
141 superintelligence.[40] Let's say that Policy X is written to do this. However, Policy X, like all other  
142 policies, is not simply just a solution to the problem, but a set of intended actions and procedures taken  
143 by the government that must first be passed by government.[41] This begs three questions: Can this  
144 policy realistically be implemented by government? How do we ensure that Policy X results in the  
145 intended outputs and outcomes? And how can we create policy and advocacy strategies to increase the  
146 chances of both of these happening? For example, while Policy X is intended to install a surveillance  
147 apparatus to prevent superintelligence, would Policy X still have that output and outcome after going  
148 through the legislature and executive branch? Is there a chance over time that it would result in  
149 mission creep? We can also develop strategies to ensure that Policy X has its intended outcomes,  
150 such as oversight mechanisms within the policy itself. We can go a step further and ask how the  
151 policymaking process itself creates implications for the AI governance field. For example, are there  
152 restrictions within the policymaking process that impacts timelines for reducing risk, such as how fast  
153 governments can act or create new laws? Could we get some form of upstream innovation from this  
154 where the policymaking process inspires or generates new ideas for AI governance?[42]

### 155 3. Typologies of AI Policy

156 Before we can talk about the policymaking process, AI policy needs to be further refined because  
157 we have to understand what kind of policies we are making. The point of this section is to show that  
158 AI risk policies are not monolithic, but rather there are multiple approaches to help achieve the same  
159 goal, and each set of these policies are going to have with it a different set of political difficulties. It  
160 also begs the question in terms of AI governance as a whole as to which sets of policies we implement  
161 and when, and which policies should be considered relevant to AI risk. In the same way that Bostrom  
162 (2014)[43] argues that there may be a preferred order of technological development, there's a similar  
163 analog with AI risk policies where there's a strategic order to policies that should be attempted to be  
164 implemented, whether it's because their political-capital cost is lower, the cost of failure is lower, or  
165 because it helps with future efforts to implement policies (such as the creation of an advisory body).

166 A typology of AI policies already has some previous explorative work to build on. Brundage  
167 (2016)[44] proposed the idea of De Facto AI policies. These are policies that already exist and are  
168 relevant to AI. These are further broken down into direct, indirect, and relevant policies. Direct  
169 policies are policies that specifically target AI such as regulations on self-driving cars. Indirect policies  
170 are policies that don't specifically target AI, but generally impact the development and diffusion of  
171 technologies (including AI) such as intellectual property laws and tort law. Relevant policies do not  
172 immediately impact AI, but are still worth considering because of their impact, such as education  
173 policy or the use of electronic medical records.

174 Brundage (2016)[45] in this paper, however, does not talk about about AI risk policy, but rather  
175 existing policies around AI as a whole. However, the classification used in this paper is overall useful  
176 and can be extended into AI risk policy. Instead of whether or not it directly or indirectly affects AI,  
177 AI risk policy can be classified into whether or not it directly or indirectly aims at reducing AI risk.  
178 Direct AI risk policies would explicitly govern the use, development, deployment, etc of AI to reduce  
179 risk. Examples of direct AI risk policy could include funding for AI safety research, rules for the  
180 development of AGI, international agreements on AI, etc. Indirect AI risk policies would either affect  
181 AI, but not explicitly govern it, or address consequences of the use of advanced AI systems. This could  
182 include both policies that affect AI and those that are AI-agnostic. For example, a policy that puts in

183 place stronger protections for privacy in general would reduce the amount of training data available,  
184 and thus the speed of AI development, and could be considered an indirect approach. An AI-agnostic  
185 policy, for example, would be basic minimum income to address technological unemployment, which  
186 could be considered a risk if it leads to societal destabilization. AI risk relevant policies would neither  
187 affect AI nor the consequences of it, but would rather make it easier for sound AI risk policies to  
188 be developed and implemented, such as changing the rules and procedures of government itself to  
189 alleviate the pacing problem.

190 There's another layer of classification that should be applied to AI risk policy based on Lowi's  
191 Typology.[46] Lowi categorizes policies into regulatory, distributive, redistributive, and constituency  
192 categories. Regulatory policies regulate one's behavior, restricting or incentivizing certain actions,  
193 such as the mandating of seat belts in cars. Distributive policies are policies that take money from the  
194 general treasury and use them for a specific project that directly benefits one group, such as a dam  
195 or research grants. Redistributive policies are those which fundamentally alter the distribution of  
196 wealth and resources in the whole of society, such as tax and welfare policies. Constituency policies  
197 are those that alter the composition and the rules and regulations of government, such as creating a  
198 new executive agency.

199 Each one of these typologies has with it a certain set of political conditions, as they impact people,  
200 businesses, and members of government differently. For example, both basic minimum income and  
201 the creation of AI safety standards are policies that are intended to reduce existential risk. However,  
202 both of these policies will have a different set of political pressures. Basic minimum income is a  
203 redistributive policy, which would move substantial amounts of wealth between classes of society.  
204 This would mean that it would likely become a nation-wide controversial issue with two opposing  
205 camps based largely on who benefits and who loses. By contrast, AI safety standards are a regulatory  
206 policy, and while there would be two groups opposed to each other on the issue (unless it comes  
207 in the form of voluntary self-regulation by the industry), the political factors around it would look  
208 different. Regulatory policies are not usually salient or popular to the general public, and thus the  
209 political battle would be largely limited to regulators, experts, and the business class. This typology  
210 will help us understand how the different policies will be treated in the policymaking process. In other  
211 words, policy creates politics. Further work on developing this might be useful for understanding the  
212 likelihood of policies being adopted and could shift strategies for which policies to pursue.

## 213 4. The Policymaking Cycle

### 214 4.1. Problem Identification, Agenda Setting, and Policy Formulation

215 The first few steps of the policymaking process—problem identification, agenda setting, and policy  
216 formulation—are usually tied together,[47] including in a so-called 'Multiple Streams Framework'.  
217 The Multiple-Streams Framework attempts to explain how policies reach the agenda when policy  
218 entrepreneurs are able to couple the policy, politics, and problems streams to open up a policy window,  
219 the opportune time when all the conditions are right to get a policy on the agenda.[48]

#### 220 4.1.1. Problem Stream

221 There are many problems in society. However, the public does not seek government intervention  
222 for many of these problems. There are some basic requirements for an issue in society to become a  
223 policy problem, which is that it is something that the public finds to be intolerable, government can  
224 do something about, and is generally seen as a legitimate area for government to work on.[49] Policy  
225 problems can also arise when there are two or more identifiable groups who enter into conflict in a  
226 policy arena for resources or positions of power.[50]

227 The first condition for an issue to be considered a policy problem is that it is something that the  
228 public or a group finds to be intolerable. Indicators such as statistics can help identify a problem.  
229 These can be used objectively, for understanding conditions in society, or politically, when they're

230 used to justify a political position: for example, using gun violence statistics as an argument for  
231 gun control. What's considered an issue over time because of the evolution of society with ever  
232 changing values, distribution of resources, technology, etc.[51] In AI governance, identifiers such  
233 as the rate of technological progress or the proliferation of autonomous weapons could be used as  
234 examples. Creating a list of politically salient identifiers or metrics could be potentially useful for  
235 creating long-term strategies and goals.

236 How the issue is framed is very important for whether or not it will be considered a policy  
237 problem.[52] Is mandating seatbelts in cars beneficial for public safety? Or is it paternalistic? Are these  
238 problems legitimate for government to handle? The framing of a problem can have an overwhelming  
239 impact on whether or not it is considered a problem appropriate for government to even formulate  
240 policy on. It can also impact the content of the policy. Whether you define access to transportation  
241 for handicapped people as a transportation problem or a civil rights issue determines whether the  
242 acceptable solution involves buying special needs vans, or costly upgrades to buses and subways to  
243 ensure equal access. Framing can also raise the priority of a policy problem by, for example, calling it a  
244 crisis and raising a sense of urgency.

245 The question of framing is also incredibly important for AI governance. For example, would  
246 autonomous weapons make war more humane by removing humans? Or will it distance ourselves  
247 from the violence and make us more willing to use them? The AI governance community needs to  
248 think about how these issues ought to be framed, and the consequences of doing so.

249 In order for an issue to be a part of the system agenda, or what the public or specific communities  
250 are discussing, there must be a focusing event. Focusing events are specific events that draw attention  
251 to a problem in society and the reasons behind it. The Sandy Hook school shooting for example is  
252 a focusing event that drew attention to America's gun laws. Moreover, events that occur outside of  
253 sector-specific focusing events,[53] or past policies on these issues, can have a large impact, especially  
254 on the types of solutions used. For AI governance, "Sputnik moments" such as AlphaGo beating Lee  
255 Sedol would be an example that drew considerable media attention and generated much discussion  
256 about the future of AI, especially in China.[54]

257 Understanding how to exploit these events for the AI governance agenda will be key to generating  
258 support and getting policies on the agenda. It's also important to stay on top of these events to  
259 understand the direction society is heading in—and to pre-empt or avert less productive or dangerous  
260 framings that might feed into arms races.[55] For example, Yampolskiy (2018) details a list of past  
261 failures by AI-enabled products.[56] How could work like this be used to influence the problem-setting?  
262 Could other AI risk researchers expand on it and build that work into a more thorough project to  
263 be used to draw attention to AI risk? Or, could attempts as this backfire and cause pre-emptive  
264 stigmatization or ineffective policies?

#### 265 4.1.2. Politics Stream

266 The politics stream is the combined factors of the national mood or public opinion, campaign  
267 groups, and administrative/legislative change. Decision-makers in government keep tabs on the  
268 swaying opinions of the masses and interest groups and act in a way that promotes themselves  
269 favorably, changing items on the agenda to stay relevant and popular, and to obscure unpopular policy  
270 stances. Changes in administration, especially when there's a major shift in the ideological composition  
271 of the institution, have a strong impact on what's included or not included on the agenda.[57]

272 In AI governance, and for people involved in advocating and implementing policies, maintaining  
273 a key eye on domestic and international politics will be key. Knowing when and what kind of policy to  
274 advocate for, and to whom, is crucial not only to saving time and energy, but also for legitimacy. Trying  
275 to sell a nationalistic administration on greater UN involvement will probably not help someone with  
276 furthering their policy proposals, and may even damage their (and their coalition's) political capital  
277 and cause. However, other forms of cooperation, such as bilateral cooperation for reducing the risk of  
278 accidents,[58] may be more promising.

279 AI governance researchers will need to consider how the political landscape should shape their  
280 recommendations or policy proposals. Not only would it determine if their recommendations would  
281 ever get considered, but if it was implemented, how would it affect the national mood? Would the next  
282 administration simply walk it back? How would other interest groups react and impact the long-term  
283 ability to reduce risk? If administration changes results in a flip-flop of ideology, what does that mean  
284 for AI risk policies associated with the past administration? Could an AI risk policy group maintain  
285 influence throughout changing administrations? All of these have implications on our ability to reduce  
286 AI risk, and this means that the policymaking strategy will not only have to be robust, but flexible  
287 enough to survive changing political conditions.

#### 288 4.1.3. Policy Stream

289 The policy stream, which is in essence the policy formulation aspect of the policy cycle, is the  
290 “soup” of ideas that are generated by policymakers,[59] when deciding what to do about a problem.  
291 Different policy networks create policies differently, with different levels of innovativeness and  
292 speed.[60] Understanding these differences, and examining their implications for the AI governance  
293 field might be useful to understand its long-term impact and the specific strategic routes it should  
294 take. In other words, how should the AI governance research field itself be organized in a way that  
295 promotes useful and relevant solutions?

296 Despite the staggering number of policy proposals coming out, only a handful will ever be  
297 accepted. These policies compete with one another and are selected on a set of criteria, which includes  
298 technical feasibility, value compatibility,[61] budgetary and political costs, and public acceptance.  
299 Policies that work will also be technically sound, with no major loopholes, and a clear rationale for  
300 how its provisions would lead to actually achieving the policy objectives.[62] This actually creates  
301 some key considerations for the field. It means that many ideas are either functionally useless due to  
302 their political limitations, unlikely to be adopted in the face of easier or less politically costly options,  
303 do not have viable policy mechanisms to achieve their goal, or are otherwise intractable prospects for  
304 government. Even if all of the above conditions are resolved, loopholes and unintended consequences  
305 may neuter the policy or make conditions worse. This vastly reduces the space of possible solutions.  
306 And even though the ability for policy implementation or values might change over time, it’s still a  
307 matter of how much and when. This begs the question: what problems can be solved when, how, and  
308 by whom? What does that mean for the large picture strategic approach?

309 Where should we seek our policies to originate from? While there are a bunch of policy ideas  
310 out there, only a few are ever seriously considered for adoption. Sources of these policies include (in  
311 the United States Federal Government, for example) the President along with the Executive Office  
312 of the President, Congressional leaders, government agencies (mostly small incremental changes  
313 and adjustments), temporary organizations or ‘ad hoc’ that serve to investigate specific topics,  
314 and interest groups whose topical expertise and political power can sometimes make them de facto  
315 policymakers. Each of these areas have differing levels of legitimacy, influence, and degree to which  
316 they can make policy changes. A question to consider is not only where in the policy network should  
317 AI risk policymakers focus on making these policies, but where they can best advocate for the creation  
318 of additional bodies like ad hoc to create additional policies, and what implications does that  
319 have for the field at large?

320 With regard to the policy formulation phase of policymaking, a continuum of political  
321 environments has been created such that on one extreme there are policies with publics and on  
322 the other, there are policies without publics.[63] When policies are formulated it is important to  
323 consider political environments relevant to the issue. The term “publics” refers to groups who have  
324 more than a passing interest in an issue or are actively involved in it. It appears that AI risks are issues  
325 where there are limited incentives for publics to form because of problems being remote, costly, or  
326 even abstract and uncertain. What does this mean for the AI safety community? How can interest

327 groups be created most effectively? How can these issues be best expressed so that they don't seem so  
328 remote, abstract, or uncertain?

#### 329 4.1.4. Policy Windows and Policy Entrepreneurs

330 This framework assumes that policy decision-makers, the legislators and bureaucrats in  
331 government, exist in a state of ambiguity, where they don't have a clear set of preferences and each set  
332 of circumstances can be seen in more than one way. This cannot be resolved with more information, as  
333 it is not an issue of ignorance. The example that Zahariadis (2007) gives is that "more information can  
334 tell us how AIDS is spread, but it still won't tell us whether AIDS is a health, educational, political, or  
335 moral issue."[\[64\]](#)

336 Overall, the Multiple Streams Framework describes government organizations as "organized  
337 anarchies" where institutional problems run rampant, there are often unclear or under-defined goals,  
338 overlapping jurisdictions, and host of other problems means that decision-makers have to ration their  
339 time between problems and do not have enough time to create a clear set of preferences, make good  
340 use of information, or take the time to comprehend the problem for sound decisions on policies. In  
341 essence, decision-makers are not rational decision-makers by any stretch. Instead, it depends on the  
342 ability of policy entrepreneurs to couple the three streams and manipulate the decision-maker into  
343 achieving their intended policy goals.[\[65\]](#)

344 Policy entrepreneurs, who are the policymakers, advocates, interest groups, etc. who push to  
345 make specific legislative changes in their areas, only have a short window of time to have their  
346 proposals added to the formal agenda. It's when the right political environment, a timely problem, and  
347 a potentially acceptable solution all meet together with a policy entrepreneur who can manipulate the  
348 situation to their advantage. Because decision-makers exist in a state of ambiguity, policy entrepreneurs  
349 are able to manipulate their interpretation of their information to provide meaning, identity, and clarity.

350 Policy entrepreneurs use different tools and tactics to manipulate the way decision-makers process  
351 information and exploit their behavioral biases. Framing tactics, for example, can be used to present  
352 a policy option as a loss to the status quo, not taking note of the degree of loss it creates, exploiting  
353 decision-makers who are loss averse, and may push them towards more extreme options like going to  
354 war to make up for those small losses.[\[66\]](#)

355 The manipulation of emotions through symbols and the identity or social status of a  
356 decision-maker can also pressure them to make certain choices; policies around flag-burning are  
357 a great example of this. Because decision-makers are under a great deal of stress and are time  
358 constrained, the strategic ordering of decisions, or 'salami tactics,' creates agreement in steps by  
359 reducing the total perceived risk of a policy.[\[67\]](#) The manipulation of symbols in the way that artificial  
360 intelligence is being framed today has already occurred. At first, anti-autonomous weapons advocates  
361 were describing 'armed quadcopters' as a serious problem with little media attention.[\[68\]](#) These were  
362 rebranded as 'slaughterbots' and a short-film was released with substantial media attention. But,  
363 what sort of long-run impact will this have on the field? While giving policymakers straight facts and  
364 solutions seems appealing, AI risk policymakers have to consider that it is impractical in reality, and  
365 may have to accept the inevitability, to policy success, of tactics like framing. Which begs the question,  
366 which tactics should they use and how? Questions like these must be considered.

367 All of this begs a serious consideration. Consider, if there are some problems that can only  
368 be resolved through state action (such as an arms race), that means that it is dependent on the  
369 policymaking process, and thus these solutions can only be passed when policy windows open. So,  
370 how many of these opportunities do AI risk policymakers get? Or, how many chances do they get  
371 to implement AI risk policies? These windows only open every once in a while and they're often in  
372 fragile conditions. For example, Bill Clinton's campaign in 1992 aimed to reform the healthcare system  
373 and made it a campaign priority, but his administration's failure to pass the bill closed the window.[\[69\]](#)  
374 In other words, what impact does this have on AI governance and policy implementation timelines  
375 and what does that mean for the field as a whole?

376 However, in order for a policy entrepreneur to manipulate decision-makers, they must have access  
377 to them, which is highly dependant on both the legitimacy of their issue, but also for the legitimacy  
378 of the group itself and their interest. One of the ways that policy entrepreneurs increase their own  
379 influence is to create new decision-points that they can exploit and to reduce access of other groups.[70]  
380 AI risk policymakers and advocates will have to find some way to gain access to decision-makers.  
381 For example, working on near-term or non-existential risk issues with AI might help someone build  
382 the social capital and network that's necessary to work on existential risks issues. This would not  
383 only make it easier people in the field to implement their solutions, but to also make themselves  
384 gatekeepers to the decision-makers, which could help with preventing policies that would increase  
385 existential risks (whether from AI or other sources) from getting through. This may be an area that  
386 needs further research. Aspects such as a group's access to decision-makers, the advocating group's  
387 legitimacy, biases of the institution,[71] and a group's ability to mobilize resources will determine what  
388 gets added to the agenda, and the AI risk community will need to work on building all of these. AI  
389 policymakers will need to develop a strategy for how to get the right people into the right places, and  
390 how to coordinate between different groups.

391 Getting on the formal agenda is a competitive process because there are fundamental limits to a  
392 decision-maker's time, and because the policy may be perceived to harm the interests of other groups.  
393 Opposing groups can use a variety of tactics such as denying that the problem exists, arguing that it is  
394 not a problem for government, or arguing that the solution would have bad societal consequences, to  
395 deny it agenda status. Other factors that could deny an issue agenda status include changing societal  
396 norms, political changes, or political leaders avoiding having to be confronted by an issue that hurts  
397 their interests. So AI policymakers will need to know how to overcome and adapt to these changing  
398 situations and other organizations preventing their policies from being adopted.

399 AI governance and policy experts will need to pay attention to the arguments being used for  
400 and against superintelligence, and whether or not this will become a political issue. Baum (2018)  
401 notes that superintelligence is particularly vulnerable to what's known as politicized skepticism,  
402 skepticism that is not based on an intellectual disagreement about the problem, based on good-faith  
403 attempts to understand the arguments, but rather to shut down concerns based out of self-interest (or  
404 a conflict of interests). Some major AI companies, and even other academics, have criticized the idea of  
405 superintelligence out of what seems to be their own self-interest as opposed to genuine concerns.[72]  
406 This would have a devastating impact on AI policy advocates in a similar way that the tobacco industry  
407 significantly impacted scientific efforts to study the public health links between tobacco and cancer.

#### 408 4.2. Policy Adoption

409 The next stage of the policy cycle is policy adoption, or when decision-makers choose an option  
410 that adopts, modifies, or abandons a policy. This does not necessarily take the form of choosing from a  
411 buffet of completed pieces of policy, but rather to take further action on a policy alternative that's more  
412 preferable and that's more likely to win approval. At this point, after much bargaining and discussion,  
413 the policy choice will only be a formality, or there will be continuous discussion and disagreement  
414 until there's a formal vote or decision made. This is an important field to analyze for AI policymakers  
415 for the obvious implication that they will want their policy proposals being chosen, and so they will  
416 need to understand and design strategies to do so. Also, as we will discuss later, when changes do  
417 occur, they can often bring with it wider changes in public policy,[73] an implication that will need to  
418 be taken into account.

419 The Advocacy Coalition Framework is a theory on policy adoption, but also incorporates every  
420 other aspect of the policy cycle with it. The theory describes the interactions of two or more 'advocacy  
421 coalitions'—that is, groups of people from a multitude of positions who coordinate together to advocate  
422 for some belief, or to implement some policy change (potentially over many fields) over an extended  
423 period of time.[74] These don't need to be a single, explicitly delineated organizations like the National  
424 Rifle Association, but could include loosely affiliated groups of organizations and/or individuals, all

425 working towards the same goal. Building and maintaining coalitions will be one of the major tasks  
426 that AI policymakers will need to work on, and so examining this framework will be highly valuable.

427 What is it that binds a coalition together? All advocacy coalitions share some form of beliefs.  
428 However, the Advocacy Coalition Framework uses a hierarchical belief system. The deepest and  
429 broadest of these are deep core beliefs, which are normative positions on human nature, hierarchy of  
430 value preferences (i.e., should we value liberty over equality?), the role of government, etc. Policy  
431 core beliefs are the next stage of the hierarchy, which involves the extension of deep core beliefs into  
432 policy areas. Both of these areas are very difficult to change, as they involve fundamental values. This  
433 actually creates an issue where, due to differing fundamental and personal values which lead to lack  
434 of interaction, different coalitions often see the same information differently, leading to distrust. Each  
435 may come to see the other side as “evil,” reducing the possibilities of cooperation and compromise.[75]

436 The deeply held convictions of what a policy subsystem ought to look like are called policy core  
437 policy preferences, and are the source of conflict between advocacy coalitions. They are the salient  
438 problems that have been the long-running issues in that area for a time. Policy core policy preferences  
439 shape the political landscape, dictating who allies with whom and who are the enemies, and what  
440 strategies coalitions take.

441 The final level of the belief hierarchy are secondary beliefs, belief that cover procedures, rules, and  
442 things of this nature. These are very narrow in scope and the easier to change, requiring less evidence  
443 and little bargaining to change.

444 Understanding the values and beliefs of different existing coalitions, groups, and individuals is key  
445 to building and maintaining new coalitions for AI policymakers. This brings up a few considerations.  
446 Since it’s difficult for conflicting coalitions to work together, will AI policymakers have to choose certain  
447 coalitions to work with? What are the costs, benefits, and the potential blowback of this? Since some  
448 policies related to AI risk are not in a mature policy field (and thus do not have established coalitions),  
449 what can be done to shape the field beforehand to their advantage and/or promote cooperation among  
450 coalitions that are likely to form? Also since secondary beliefs are relatively easy to change, what can  
451 be changed to help reduce existential risk?

452 On a macro-level, this AC Framework acts as a cycle. Relatively stable parameters, as mentioned  
453 before, exist in the status quo since policy arenas usually come to some equilibrium where one coalition  
454 dominates the policy subsystem. Then, policy changes made by an advocacy coalition or an outside  
455 event creates a fundamental change in the world, whether its a change in public opinion or in the  
456 rules and procedures governing a subsystem, which changes the initial stable parameters, such as a  
457 major event like a mass shooting. These lead to a shift in power that allows another coalition to gain  
458 influence over the types of policies being adopted. However, especially in the case of controversial  
459 legislature, policies that require multiple veto points to pass will create access for multiple coalitions.  
460 This means that even a coalition that dominates a subsystem won’t have unilateral ability to dictate  
461 policies in some situations. Others however, especially where there are few decision-makers or an  
462 exceptionally influential decision-maker, can result in highly monopolized systems. Questions such as  
463 how to be resilient to these changes in conditions, how to facilitate changes into conditions that are  
464 beneficial to AI policymakers, and how to construct policy subsystems in a way that’s conducive to AI  
465 policymakers’ goals are useful questions to consider.

466 This theory describes policy adoption on a very broad level, but how do the decision-makers  
467 themselves decide which policies to move forward with? Different incentives and restrictions come to  
468 play at different levels of policymaking. For example, highly salient and popular issues are more likely  
469 to be influenced by popular opinion, whereas obscure technical issues will likely be determined by  
470 policy experts in that field. Different factors that affect both individual and group decision-makers  
471 also come into play, such as their personal, professional, organizational, and ideological values.  
472 For legislators, their political party and their constituency also play an overwhelming role in their  
473 decision-making. Understanding and mapping out these factors will be necessary for the successful  
474 implementation of AI risk policy.

475 On top of these factors, decision-makers usually never have the time, expertise, or even care  
476 enough to be able to come up with a fully rational approach to deciding most policies. In many cases,  
477 legislators will seek out the advice of other legislators and experts and follow their lead. Due to this  
478 being a widespread practice, a few key institutions and leaders often have disproportionate power. For  
479 those working in AI risk policy, it's necessary to understand these things so that the message they craft  
480 for as to why policy change should occur, and who to specifically target to get widespread adoption  
481 from other decision-makers in the policy arena.

#### 482 4.3. Policy Implementation

483 Policy implementation is one step in the policymaking process. It is defined as whatever is done  
484 to carry a law into effect, to apply it to the target population. . . and to achieve its goals.“[76] In other  
485 words, it's the activity where adopted policies are carried into effect.[77] However, it is not to say that  
486 it's a very distinct step that can be clearly distinguished from others. Every implementation action can  
487 influence policy problems, resources, and objectives as the process evolves.[78] Policy implementation  
488 can influence problem identification, policy adoption, etc.

489 Two broad factors that have been offered for the success of policy are local capacity and will.[79]  
490 In other words, is there enough training, money, and human resources, along with the right attitudes,  
491 motivation, and beliefs to make something happen? It is suggested that the former can be influenced  
492 much more easily than the latter as more money can be received and consultants can be hired. For  
493 AI risk, both questions are relevant: how to increase capacity and how to influence the influencers.  
494 With the former, it has been estimated that about \$9[80]-\$20[81] million is currently spent on AI risk.  
495 With the latter, studying the opinion of the public as well as experts might be a useful approach. One  
496 survey[82] indicates that only 8% of top-cited authors in AI consider that human-level AI would be  
497 extremely bad (existential risk) for humanity. Another survey that is more recent[83] indicates that  
498 machine learning researchers think on average (median) that there's a 10% probability that human-level  
499 machine intelligence will result in a negative outcome and 5% probability that it will have an extremely  
500 bad outcome (existential risk). The general public seems to be generally cautious, with a survey  
501 showing 82% of Americans believing that AI/robots should be managed carefully.[84]

502 This part of the policymaking process is very difficult as the literature is generally quite pessimistic  
503 about the ability of policies to bringing social changes into effect.[85] However, the authors of the cited  
504 paper have identified conditions of effective implementation based on successful examples. These  
505 conditions are a) the policy is based on a sound theory of getting the target group to behave in a  
506 desired way, b) policy directives and structures for the target group are unambiguous, c) the leaders  
507 implementing the policies are skillful with regard to management and politics and committed to the  
508 goals, d) policy is supported by organized constituency groups and key legislators as well as courts  
509 throughout the implementation process, and e) the relative priority of policies is not significantly  
510 undermined over time by other policies or socioeconomic changes. Additionally,[86] having carefully  
511 drafted statute that incentivizes behavior changes, provides adequate funds, expresses clearly ranked  
512 goals, is an implementation process, and has few veto points, is also vital to the success of a policy.

513 In regards to AI governance, the ambiguity and complexity of the problem creates a major hurdle  
514 for effective policies to be developed. Breaking down AI risk policy into multiple domains as discussed  
515 in the previous section helps with creating somewhat less ambiguous objectives, such as changing the  
516 education system to be more conducive for technological growth. Even then however, because many  
517 of the issues are either complex or have not happened yet, it's difficult to create concrete objectives  
518 and policies. AI risk is not like noise pollution, where there's an easily identifiable, manageable, and  
519 tractable problem. Further research could help identify concrete and tractable issues that might lead to  
520 a reduction of risk. In addition, when trying to develop and implement policy, AI policymakers will  
521 need to keep in mind factors such as what extent there is support for it in the executive branch, with  
522 outside organizations, and how exactly the policy is written and how those change throughout the  
523 policymaking cycle.

524 Another key consideration for successful policy implementation that we were able to identify  
525 from the literature is engaging with the community to increase readiness to accept and devote resources  
526 to policy-related problems. It has been acknowledged that there are no good evidence-based ways  
527 of achieving community buy-in. This is an area that might be useful to study in order to increase the  
528 chances of successful reduction of AI risk. There are different stages of community readiness such as  
529 no awareness, denial, and vague awareness to preplanning, preparation, initiation, and stabilization  
530 phases.[87] It's important to understand what counts as the community and what phases different  
531 subcommunities of AI safety field are in. Earlier, we mentioned a survey about AI experts and showed  
532 that their readiness with AI risks was low. Other relevant experts, the public, and other types of  
533 subcommunities might have different levels of readiness.

534 It has been suggested that "the more clearly the core components of an intervention program  
535 or practice are known and defined, the more readily the program or practice can be implemented  
536 successfully." [88] In other words, policies and steps of implementation of those policies have to be  
537 very clearly expressed. What implications does this have for AI risk? Researchers and policymakers  
538 should evaluate how clearly core components have been expressed in this field and improve them as  
539 necessary.

#### 540 *4.4. Policy Evaluation*

541 The final step in the policymaking cycle is policy evaluation. This includes activities related  
542 to determining the impact of the policy, whether it is achieving its goals, whether the rules and  
543 procedures it lays out are being followed, and other externalities or unintended consequences.[89] As  
544 we've explained before, policy evaluation doesn't have to occur only at this step. For example, the  
545 impact of a policy is estimated already in the early stages. Anderson et al. highlight different types of  
546 policy evaluations in their book, but especially consider systematic evaluations of programs. It involves  
547 "the specification of goals or objectives; the collection of information and data on program inputs,  
548 outputs, and consequences; and their rigorous analysis, preferably through the use of quantitative or  
549 statistical techniques." [90]

550 Policy evaluation examines a policy to understand its impacts in multiple ways.[91] First, is the  
551 policy affecting the population that it's intending to target? In AI risk policy, this could be anything  
552 from large tech companies, to AI researchers, to people affected by technological unemployment.  
553 Second, are there populations that are being affected that were not intended? These externalities  
554 could be positive or negative. Third, what are the benefits and costs associated with this policy? AI  
555 policymakers will want to ensure that their policies actually reduce risk and that the costs are not so  
556 astronomical that they become politically infeasible. Finally, what long-term costs and benefits does  
557 a policy have? This is especially important for AI risk policy, as decisions now could have a major  
558 impact on the long-term risk that AI has. In AI governance and policymaking, research needs to be  
559 done on what sort of indicators or metrics are used for the reduction of risk, and for identifying what  
560 goals that should be achieved.

561 If the previous steps in the policymaking process have generated goals that are unclear or diverse,  
562 it's very difficult to evaluate the impact of the policy.[92] Different decision-makers can more easily  
563 reach differing conclusion about the results of a program in that case, or may not follow it all all.[93]  
564 How the goals of an AI risk program are defined is, therefore, very important.

565 Another key consideration for policy evaluation is how to make sure that the results are objectively  
566 measured. Agency and program officials may be wary of possible political consequences of the  
567 evaluation process.[94] If it turns out that the program was not useful or even detrimental, this might  
568 have consequences to their influence and career. Because of this consideration, they might not be very  
569 interested in correct evaluation studies or they may hinder the process in some other way. There are  
570 many ways an evaluation of a policy might be ignored or attacked such as claiming it was poorly done,  
571 the data was inadequate, or the findings inconclusive.[95] Thus, it is important that researchers are  
572 provided with high quality and relevant data-sets that

573 There is also the distinction between policy outputs and outcomes[96] to consider. Outputs are  
574 tangible actions taken or things produced, such as collecting taxes or building a dam. Outcomes on the  
575 other hand, are the consequences for society, such as lower disposable income or cleaner air quality.  
576 Outputs do not always produce the intended outcomes, which is highly evident in areas such as social  
577 welfare policy, where policies may unintentionally trap people in poverty. For AI policymakers, this is  
578 very important to consider whether their policy outputs will have the intended consequences, and if  
579 so, how to correct that policy.

580 The evaluation of a policy and the political responses to it can result in the termination of  
581 it.[97] Assuming that AI risk policymakers do not want their policies to be terminated or altered in  
582 a detrimental way, how can they make sure it doesn't happen? A policy getting altered to be more  
583 effective might be a good thing, but termination can bring unpleasant and negative connotations. It  
584 might even have negative consequences to the community.[? ] What exact consequences might it have  
585 politically? Also, it's important to remember that many policymaker's time horizon only goes till the  
586 next election, and so they often seek immediate results, often before the returns come into fruition.  
587 While this may not impact all policies, as this mostly applies to salient policies like healthcare and  
588 education, AI policymakers should keep this in mind and try to understand how it might impact their  
589 work.

## 590 5. Conclusion

591 There are multiple policy options that could be chosen that either directly or indirectly reduce AI  
592 risk, or relevant policies that could help with further efforts to reduce AI risk. Because different policy  
593 arenas have different political conditions, and the policymaking process itself draws a number of  
594 important challenges, it begs the questions as to what policies in what order are chosen, what strategies  
595 are used to get these policies passed and implemented by the government, and the larger impact of  
596 these choices on AI governance and risk as a whole. We argued that a new subfield of AI governance  
597 research on AI policymaking strategies should be further investigated to draw implications for how  
598 these policies should be designed, advocated for, and how organizations should approach solving this  
599 issue.

## 600 6. Limitations and Future Research

601 This paper is intended to be a broad overview to be a conversation starter for future research  
602 into this area. Thus, there is a strong limitation to the depth of research in this paper. However, we  
603 expect that future work will be done to further refine the line of thinking we laid out above, along with  
604 further in-depth study into the different theories and their applicability to AI risk.

605 One of the major limitations of this paper is that the stages heuristic presented in this paper has  
606 been heavily criticized and is subject to debate about its effectiveness. Sabatier (2007) has criticized it  
607 for not being a causal theory, having a strong top-down bias, among other critiques. However, he also  
608 notes that there is much up to debate, with some scholars such as Anderson (2010) advocating for it.  
609 There are also a number of other theories that were not discussed in this paper, such as Institutional  
610 Rational Choice, the Punctuated-Equilibrium Framework, the Policy Diffusion Framework, and other  
611 lesser known theories. We expect that future research will explore which policy frameworks should be  
612 focused on in AI risk research.

613 The other limitation of this paper is that its applicability to the international governance of AI  
614 was not discussed. Future research that looks at how much these theories apply to foreign policy and  
615 the international governance of AI in general would be useful. If these theories have a very limited or  
616 no impact on the international governance of AI, then figuring out how much work can be done to  
617 reduce AI risk in domestic policy would determine the usefulness of these theories.

618 Throughout the paper, we have come up with a number of key considerations. For convenience,  
619 we have come up with a list of these considerations below.

## 620 7. Appendix/Summary

621 In this part, we summarize key considerations we were able to find and come up with when we  
622 reviewed literature on policymaking process.

623 Thesis level consideration:

- 624 • What problems can be solved when, how, and by whom? What does that mean for the large  
625 picture strategic approach?

626 Considerations from Typologies of Policies:

- 627 • Are there AI risk policies that should be implemented first? How would we decide this?
- 628 • What types of policies should the AI risk policymakers try to get implemented? Why should  
629 those types be prioritized?
- 630 • What the political considerations surrounding different sets of policies and how does that affect  
631 its ability to be implemented?

632 Considerations from Problem Identification, Agenda Setting, and Policy Formulation:

- 633 • Is this issue or policy legitimated?
- 634 • Would the policy be supported by the current administration and be able to be maintained  
635 through changing administrations?
- 636 • Which policies out of different sets of potential solutions are politically feasible?
- 637 • Are there less costly alternative policies that AI risk policymakers will have to compete with?
- 638 • How does attention to problems by different communities affect AI risk policymaker's actions?
- 639 • What types of framing of policy issues is most beneficial? What types are most dangerous?
- 640 • Is there a way to determine how framing will determine policy content?
- 641 • What focusing events have occurred in the field of AI?
- 642 • How can AI risk policymakers utilize focusing events to further policy agendas?
- 643 • What effect do other organizations have on reducing the legitimacy of AI risk?
- 644 • What can be done to respond to these counter-movements effectively? What kind of responses to  
645 objections are most convincing?
- 646 • How many policy windows will there be for a particular issue? What does this mean for AI risk  
647 policymaker's overall strategy?
- 648 • What role should AI risk policy entrepreneurs play in AI governance?
- 649 • How and where should AI risk policy entrepreneurs gain access in government?

650 Considerations from Policy Adoption:

- 651 • What policy alternatives are more likely to win approval to improve the odds of success for AI  
652 risk reduction?
- 653 • What strategies can be used to improve the chances of a preferred policy to be adopted?
- 654 • Which groups or individuals could join AI risk coalitions, what criteria are used to decide this,  
655 and what costs does them joining the coalition have?
- 656 • What role can organizations outside of AI risk play in furthering AI risk policymaker's agenda?

657 Considerations from Policy Implementation:

- 658 • Is this solution technically feasible for governments to implement?
- 659 • Are there enough resources, will, and support by leaders and constituency groups to be successful  
660 in implementation?
- 661 • Is the policy crafted in a way that effectively structures incentives for the target group?
- 662 • Is the policy unambiguous? If so, then how will that affect its ability to be implemented?
- 663 • Are the goals of the policy in conflict with any other policy or changes in society?
- 664 • Are there any veto points in the policy's statutes to prevent effective implementation?

- 665 • How will the contents or the political factors surrounding of a policy be affected during
- 666 implementation?
- 667 • Do the relevant communities accept the issue and are they willing to devote resources to resolve
- 668 it?

#### 669 Considerations from Policy Evaluation:

- 670 • Are the policy outputs having the intended outcomes?
- 671 • What are the consequences of any unintentional outcomes?
- 672 • What are the political factors surrounding the metrics that are being used to evaluate the policy?
- 673 • Do the political costs or benefits of the policy have an impact on its success?
- 674 • If the policy is terminated, will there be any negative political consequences?
- 675 • How can AI risk policymakers update the policy? How can they prevent changes by other
- 676 groups that would be harmful?
- 677 • How will the limited time horizons of lawmakers and other groups affect the evaluation of the
- 678 policy?

#### 679 References

- 680 1. Tegmark, Max. *Life 3.0: Being Human in the Age of Artificial Intelligence*. New York: Knopf, 2017.
- 681 2. Bostrom, N. *Superintelligence: Paths, Dangers, Strategies*; Oxford University Press: Oxford, New York, 2014.
- 682 3. Dafoe, Allan. "AI Governance: A Research Agenda." Oxford: Governance of AI Program, Future of
- 683 Humanity Institute, 2018. <https://www.fhi.ox.ac.uk/govaiagenda/>.
- 684 4. Baum, S. D. A Survey of Artificial General Intelligence Projects for Ethics, Risk, and Policy. Global
- 685 Catastrophic Risk Institute Working Paper 17-1. 2017. Available online: [https://papers.ssrn.com/sol3/](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3070741)
- 686 [papers.cfm?abstract\\_id=3070741](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3070741) (accessed on 11 November 2019).
- 687 5. Everitt, T.; Lea, G.; Hutter, M. AGI Safety Literature Review. ArXiv180501109 Cs 2018
- 688 6. Joy, B. Why the future doesn't need us. *Wired*, 8(4), 238-263. 2000. Available online: [https://www.wired.](https://www.wired.com/2000/04/joy-2/)
- 689 [com/2000/04/joy-2/](https://www.wired.com/2000/04/joy-2/) (accessed on 06 January 2019).
- 690 7. Hibbard, B. *Super-Intelligent Machines*. New York: Springer, 2002.
- 691 8. Hughes, J.J. Global technology regulation and potentially apocalyptic technological threats. In *Nanoethics:*
- 692 *The Ethical and Social Implications of Nanotechnology*; Allhoff, F., Eds.; Hoboken, NJ: John Wiley, 2007; pp.
- 693 201-214.
- 694 9. McGinnis, John O. "Accelerating AI." *Northwestern University Law Review*. 2010. Pp 104. Available
- 695 online: [https://scholarlycommons.law.northwestern.edu/cgi/viewcontent.cgi?referer=&httpsredir=1&](https://scholarlycommons.law.northwestern.edu/cgi/viewcontent.cgi?referer=&httpsredir=1&article=1193&context=nulr_online)
- 696 [article=1193&context=nulr\\_online](https://scholarlycommons.law.northwestern.edu/cgi/viewcontent.cgi?referer=&httpsredir=1&article=1193&context=nulr_online). (Accessed 14 March 2019)
- 697 10. Everitt, T.; Lea, G.; Hutter, M. AGI Safety Literature Review. ArXiv180501109 Cs 2018.
- 698 11. Scherer, M.U. Regulating artificial intelligence systems: risks, challenges, competencies, and strategies. *Harv.*
- 699 *J. Law Technol.* 2016, 29, 353(48).
- 700 12. Guihot, M.; Matthew, A.F.; Suzor, N.P. Nudging robots: Innovative solutions to regulate artificial intelligence.
- 701 *Vanderbilt J. Entertain. Technol. Law* 2017, 20, 385–456.
- 702 13. Baum, S.D. On the promotion of safe and socially beneficial artificial intelligence. *AI Soc.* 2017, 32, 543–551.
- 703 14. Yampolskiy, R.; Fox, J. Safety Engineering for Artificial General Intelligence. *Topoi* 2013, 32, 217–226.
- 704 15. Erdelyi, Olivia J., and Judy Goldsmith. "Regulating Artificial Intelligence: Proposal for a Global Solution." AAI/ACM Conference on Artificial Intelligence Ethics, and Society. New Orleans, USA: Association
- 705 for the Advancement of Artificial Intelligence. Social Science Research Network. 2018. Available online:
- 706 [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3263992](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3263992) (accessed on 06 January 2019).
- 707
- 708 16. Wilson, G. Minimizing global catastrophic and existential risks from emerging technologies through
- 709 international law. *Va. Environ. Law J.* 2013, 31, 307–364.
- 710 17. Bostrom, N. *Superintelligence: Paths, Dangers, Strategies*; Oxford University Press: Oxford, New York, 2014.
- 711 18. Shulman, C. Arms control and intelligence explosions. In *7th European Conference on Computing and*
- 712 *Philosophy (ECAP)*, Bellaterra, Spain, 2-4 July 2009.
- 713 19. Goertzel, B. The Corporatization of AI is a Major Threat to Humanity. *h+ Magazine*. 2017. Available
- 714 online: <http://hplusmagazine.com/2017/07/21/corporatization-ai-major-threat-humanity/> (accessed on
- 715 06 January 2019).

- 716 20. Goertzel, B. The Corporatization of AI is a Major Threat to Humanity. *h+ Magazine*. 2017. Available  
717 online: <http://hplussmagazine.com/2017/07/21/corporatization-ai-major-threat-humanity/> (accessed on  
718 06 January 2019).
- 719 21. Bostrom, N. Strategic Implications of Openness in AI Development. *Glob. Policy* 2017, 8, 135–148.
- 720 22. Shulman, C. Arms control and intelligence explosions. In 7th European Conference on Computing and  
721 Philosophy (ECAP), Bellaterra, Spain, 2-4 July 2009.
- 722 23. Dewey, D. Long-term strategies for ending existential risk from fast takeoff. In *Risks of Artificial Intelligence*;  
723 Müller V.C., Eds.; Boca Raton: CRC, 2015; pp. 243-266.
- 724 24. Goertzel, B. Should Humanity Build a Global AI Nanny to Delay the Singularity Until It's Better Understood?  
725 *J. Conscious. Stud.* 2012, 19, 96.
- 726 25. Dafoe, A. AI Governance: A Research Agenda. 2018. Available online:  
727 <https://www.fhi.ox.ac.uk/wp-content/uploads/GovAIAgenda.pdf> (accessed on 17 December 2018).
- 728 26. Brundage, M. et al. "The Malicious Use of Artificial Intelligence: Forecasting, Prevention,  
729 and Mitigation." Available online: [https://img1.wsimg.com/blobby/go/3d82daa4-97fe-4096-9c6b-  
730 376b92c619de/downloads/1c6q2kc4v\\_50335.pdf](https://img1.wsimg.com/blobby/go/3d82daa4-97fe-4096-9c6b-376b92c619de/downloads/1c6q2kc4v_50335.pdf) (accessed on 06 January 2018).
- 731 27. Whittlestone, Jess, Rune Nyrup, Anna Alexandrova, and Stephen Cave. "The Role and Limits of Principles  
732 in AI Ethics: Towards a Focus on Tensions." In *Proceedings of AAAI / ACM Conference on Artificial  
733 Intelligence, Ethics and Society 2019*, 7, 2019. [http://www.aies-conference.com/wp-content/papers/main/  
734 AIES-19\\_paper\\_188.pdf](http://www.aies-conference.com/wp-content/papers/main/AIES-19_paper_188.pdf).
- 735 28. Calo, R. Artificial Intelligence Policy: A Primer and Roadmap. 2017. Available online: [https://ssrn.com/  
736 abstract=3015350](https://ssrn.com/abstract=3015350) (accessed 06 January 2019). It should also be noted that Calo is dismissive of the risk of  
737 artificial general intelligence.
- 738 29. Everitt, T.; Lea, G.; Hutter, M. AGI Safety Literature Review. *ArXiv180501109 Cs* 2018.
- 739 30. Calo, R. Artificial Intelligence Policy: A Primer and Roadmap. 2017. Available online: [https://ssrn.com/  
740 abstract=3015350](https://ssrn.com/abstract=3015350) (accessed 06 January 2019).
- 741 31. Scherer, M. "Regulating Artificial Intelligence Systems: Risks, Challenges, Competencies, and Strategies." In  
742 *Harvard Journal of Law & Technology* Vol 29, Num 2. 2016.
- 743 32. Everitt, T.; Lea, G.; Hutter, M. AGI Safety Literature Review. *ArXiv180501109 Cs* 2018.
- 744 33. See also Cave, Stephen, and Seán S. Ó hÉigeartaigh. "An AI Race for Strategic Advantage: Rhetoric  
745 and Risks." In *AAAI / ACM Conference on Artificial Intelligence, Ethics and Society, 2018*. Available  
746 online: [http://www.aies-conference.com/wp-content/papers/main/AIES\\_2018\\_paper\\_163.pdf](http://www.aies-conference.com/wp-content/papers/main/AIES_2018_paper_163.pdf). (Accessed  
747 14 March 2019)
- 748 34. Bostrom, N. *Superintelligence: Paths, Dangers, Strategies*; Oxford University Press: Oxford, New York, 2014.
- 749 35. Dafoe, A. AI Governance: A Research Agenda. 2018. Available online: [https://www.fhi.ox.ac.uk/wp-  
750 content/uploads/GovAIAgenda.pdf](https://www.fhi.ox.ac.uk/wp-content/uploads/GovAIAgenda.pdf) (accessed on 17 December 2018).
- 751 36. *Ibid.*
- 752 37. Flynn, C. Personal thoughts on careers in AI policy and strategy. *Effective Altruism Forum*. 2017. Available  
753 online: [https://forum.effectivealtruism.org/posts/RCvetzfDnBNFX7pLH/  
754 personal-thoughts-on-careers-  
in-ai-policy-and-strategy](https://forum.effectivealtruism.org/posts/RCvetzfDnBNFX7pLH/personal-thoughts-on-careers-in-ai-policy-and-strategy) (accessed on 06 January 2019).
- 755 38. Which seems to have been later redefined to 'AI governance' in the Future of Humanity Institute's AI  
756 Governance: A Research Agenda
- 757 39. Goertzel, B. Should Humanity Build a Global AI Nanny to Delay the Singularity Until It's Better Understood?  
758 *J. Conscious. Stud.* 2012, 19, 96.
- 759 40. The specifics issues will depend on the type of government. For example, the types of difficulties would be  
760 different in a democracy vs a dictatorship. This paper however will focus on federal republics.
- 761 41. Thank you to Sabrina Kavanagh for suggesting the idea that the policy process could inspire new ideas for  
762 AI governance researchers.
- 763 42. Bostrom, N. *Superintelligence: Paths, Dangers, Strategies*; Oxford University Press: Oxford, New York, 2014.
- 764 43. Brundage, M.; Bryson, J. *Smart Policies for Artificial Intelligence*. *ArXiv160808196 Cs* 2016.
- 765 44. *Ibid.*
- 766 45. Lowi, T.J. Four Systems of Policy, Politics, and Choice. *Public Adm. Rev.* 1972, 32, 298–310.
- 767 46. Anderson, J.E. *Public Policymaking: An Introduction*; 7 edition.; Cengage Learning: Boston, MA, 2010. pp 5.

- 768 47. Zahariadis, N. The Multiple Streams Framework: Structure, Limitations, Prospects. In *Theories of the Policy*  
769 *Process* 2nd Edition; Sabatier, P.; Eds.; Boulder: Westview Press, 2007; pp 70-74
- 770 48. Anderson, J.E. *Public Policymaking: An Introduction*; 7 edition.; Cengage Learning: Boston, MA, 2010. pp  
771 85-88.
- 772 49. Cobb, Roger and Charles D. Elder. What is an Issue? What Makes an Issue? In *Participation in American*  
773 *Politics: the Dynamics of Agenda Building* 1983, 82-93. Baltimore: Johns Hopkins University Press.
- 774 50. Anderson, J.E. *Public Policymaking: An Introduction*; 7 edition.; Cengage Learning: Boston, MA, 2010. pp  
775 85-88.
- 776 51. Ibid, 88.
- 777 52. Zahariadis, N. The Multiple Streams Framework: Structure, Limitations, Prospects. In *Theories of the Policy*  
778 *Process* 2nd Edition; Sabatier, P.; Eds.; Boulder: Westview Press, 2007; pp 84.
- 779 53. Allen, G. China's Artificial Intelligence Strategy Poses a Credible Threat to U.S. Tech Leadership. Center for  
780 Foreign Affairs Blog. Available online: [https://www.cfr.org/blog/chinas-artificial-intelligence-strategy-](https://www.cfr.org/blog/chinas-artificial-intelligence-strategy-poses-credible-threat-us-tech-leadership)  
781 [poses-credible-threat-us-tech-leadership](https://www.cfr.org/blog/chinas-artificial-intelligence-strategy-poses-credible-threat-us-tech-leadership) (accessed on February 26, 2019).
- 782 54. Yampolskiy, R. Current State of Knowledge on Failures of AI Enabled Products. Report. Consortium for Safer  
783 AI. 2018. Available online: [https://docs.wixstatic.com/ugd/ace275\\_0ea60fe9b665439bb0b37d20beb89b6f.](https://docs.wixstatic.com/ugd/ace275_0ea60fe9b665439bb0b37d20beb89b6f.pdf)  
784 [pdf](https://docs.wixstatic.com/ugd/ace275_0ea60fe9b665439bb0b37d20beb89b6f.pdf) (accessed on 06 January 2018).
- 785 55. Zahariadis, N. The Multiple Streams Framework: Structure, Limitations, Prospects. In *Theories of the Policy*  
786 *Process* 2nd Edition; Sabatier, P.; Eds.; Boulder: Westview Press, 2007; pp 73.
- 787 56. Danzig, R. Technology Roulette: Managing Loss of Control as Many Militaries Pursue Technological  
788 Superiority. Center for New American Security. May 2018. Available online: [https://www.cnas.org/](https://www.cnas.org/publications/reports/technology-roulette)  
789 [publications/reports/technology-roulette](https://www.cnas.org/publications/reports/technology-roulette) (Accessed on 24 March 2019).
- 790 57. Ibid, 72.
- 791 58. Ibid.
- 792 59. Ibid.
- 793 60. Anderson, J.E. *Public Policymaking: An Introduction*; 7 edition.; Cengage Learning: Boston, MA, 2010. pp  
794 108.
- 795 61. May, P.J. Reconsidering Policy Design: Policies and Publics. *J. Public Policy* 1991, 11, 187–206.
- 796 62. Zahariadis, N. The Multiple Streams Framework: Structure, Limitations, Prospects. In *Theories of the Policy*  
797 *Process* 2nd Edition; Sabatier, P.; Eds.; Boulder: Westview Press, 2007; pp 66.
- 798 63. Ibid, 74-78.
- 799 64. Ibid, 76.
- 800 65. Ibid, 65-92.
- 801 66. Russell, S.; Aguirre, A.; Conn, A.; Tegmark, M. Why You Should Fear "Slaughterbots" - A Response. *IEEE*  
802 *Spectrum*. 2018. Available online: [https://spectrum.ieee.org/automaton/robotics/artificial-intelligence/](https://spectrum.ieee.org/automaton/robotics/artificial-intelligence/why-you-should-fear-slaughterbots-a-response)  
803 [why-you-should-fear-slaughterbots-a-response](https://spectrum.ieee.org/automaton/robotics/artificial-intelligence/why-you-should-fear-slaughterbots-a-response) (accessed on 09 January 2019).
- 804 67. Zahariadis, N. The Multiple Streams Framework: Structure, Limitations, Prospects. In *Theories of the Policy*  
805 *Process* 2nd Edition; Sabatier, P.; Eds.; Boulder: Westview Press, 2007; pp 70-74
- 806 68. Cobb, R.; Elder, C.D. What is an Issue? What Makes an Issue? In *Participation in American Politics: the*  
807 *Dynamics of Agenda-Building*; Baltimore: Johns Hopkins University Press, 1983; pp. 82-93.
- 808 69. Yudkowsky, E. Cognitive Biases Potentially Affecting Judgment of Global Risks. In *Global Catastrophic*  
809 *Risks*, edited by Nick Bostrom and Milan M. Ćirković, New York: Oxford University Press. 2008. pp 91–119.
- 810 70. Baum, S.D. Superintelligence Skepticism as a Political Tool. *Information* 2018, 9, 209.
- 811 71. James, T.L.; Jones B.D.; Baumgartner F.R.. Punctuated-Equilibrium Theory: Explaining Stability and Change  
812 in Public Policymaking. In *Theories of the Policy Process*, 2nd ed. Sabatier, P.A., ed. Chapter 6. Boulder,  
813 Colorado: Westview Press. 2007
- 814 72. Sabatier, P.; Weiblle, C.M. An Advocacy Coalition Framework. In *Theories of the Policy Process*, 2nd ed.  
815 Sabatier, P.A., ed. Chapter 7. Boulder, Colorado: Westview Press. 2007.
- 816 73. Sabatier, P.; Weiblle, C.M. An Advocacy Coalition Framework. In *Theories of the Policy Process*, 2nd ed.  
817 Sabatier, P.A., ed. Chapter 7. Boulder, Colorado: Westview Press. 2007.
- 818 74. Anderson, J.E. *Public Policymaking: An Introduction*; 7 edition.; Cengage Learning: Boston, MA, 2010. pp  
819 209.
- 820 75. Ibid.

- 821 76. McLaughlin, M.W. Learning From Experience: Lessons From Policy Implementation. *Educ. Eval. Policy*  
822 *Anal.* 1987, 9, 171–178.
- 823 77. Ibid.
- 824 78. Farquhar, S. Changes in funding in the AI safety field. 2017. Available online: [https://www.  
825 centreforeffectivealtruism.org/blog/changes-in-funding-in-the-ai-safety-field](https://www.centreforeffectivealtruism.org/blog/changes-in-funding-in-the-ai-safety-field) (accessed online 06 January  
826 2019).
- 827 79. MacAskill, W. What are the most important moral problems of our time? TED Talk. 2018. Available online:  
828 [https://www.ted.com/talks/will\\_macaskill\\_how\\_can\\_we\\_do\\_the\\_most\\_good\\_for\\_the\\_world](https://www.ted.com/talks/will_macaskill_how_can_we_do_the_most_good_for_the_world) (accessed  
829 online 06 January 2019).
- 830 80. Müller, V; Bostrom, N. Future progress in artificial intelligence: A Survey of Expert Opinion, in Vincent C.  
831 Müller (ed.), *Fundamental Issues of Artificial Intelligence* (Synthese Library; Berlin: Springer). (forthcoming  
832 2014) Available online: <https://nickbostrom.com/papers/survey.pdf> (accessed online 06 January 2019).
- 833 81. Grace, K.; Salvatier, J.; Dafoe, A.; Zhang, B.; Evans, O. When Will AI Exceed Human Performance? Evidence  
834 from AI Experts. *ArXiv170508807 Cs* 2017.
- 835 82. Zhang, B.; Dafoe, A. Artificial Intelligence: American Attitudes and Trends. January 2019. Available  
836 online: [https://governanceai.github.io/US-Public-Opinion-Report-Jan-2019/us\\_public\\_opinion\\_report\\_  
837 jan\\_2019.pdf](https://governanceai.github.io/US-Public-Opinion-Report-Jan-2019/us_public_opinion_report_jan_2019.pdf) (accessed online 03 January 2019).
- 838 83. Sabatier, P.; Mazmanian, D. The Conditions of Effective Implementation: A Guide to Accomplishing Policy  
839 Objectives. *Policy Anal.* 1979, 5, 481–504.
- 840 84. Sabatier, P.; Mazmanian, D. The Implementation of Public Policy: A Framework of Analysis\*. *Policy Stud. J.*  
841 1980, 8, 538–560.
- 842 85. Ibid, 10.
- 843 86. Ibid, 24.
- 844 87. Anderson, J.E. *Public Policymaking: An Introduction*; 7 edition.; Cengage Learning: Boston, MA, 2010. pp  
845 245.
- 846 88. Ibid, 246-247.
- 847 89. Ibid, 246.
- 848 90. Ibid, 260-261.
- 849 91. Ibid, 246.
- 850 92. Ibid, 263.
- 851 93. Ibid, 264.
- 852 94. Ibid, 248-249.
- 853 95. Anderson, J.E. *Public Policymaking: An Introduction*; 7 edition.; Cengage Learning: Boston, MA, 2010. pp  
854 273.
- 855 96. Ibid, 275.
- 856 97. Ibid, 276.

857 © 2019 by the authors. Submitted to *Big Data Cogn. Comput.* for possible open access  
858 publication under the terms and conditions of the Creative Commons Attribution (CC BY) license  
859 (<http://creativecommons.org/licenses/by/4.0/>).