

Proposing a Model of Computational Responsibility

William T. Wallis

School of Computing Science Sir Alwyn Williams Building University of Glasgow G12 8QQ

Masters project proposal

Contents

1	Intr	Introduction		
	1.1	An ea	rly rebuttal of some common criticisms	3
	1.2	Termi	nology	4
2	Stat	ement (of Problem	5
	2.1	Reflec	tive Agents	6
	2.2	Interp	retive Agents	7
	2.3	Reacti	ve Agents	8
	2.4	Using	the definitions	9
3 Background Surv		kgroun	d Survey	10
	3.1	Social	Sciences and Mathematics	10
		3.1.1	Birkhoff's Aesthetic Measure	10
		3.1.2	Deutsch	11
		3.1.3	Luhmann	12
	3.2	Mode	rn [Computational] Trust methods	13
		3.2.1	Marsh's formalism	13
		3.2.2	Castelfranchi & Falcone	14
		3.2.3	Ian Sommerville, Sociotechnical Systems, and Responsibility Modelling	14
	3.3	Philos	ophy of Moral Responsibility	16
		3.3.1	Peter F. Strawson	16
		3.3.2	Thomas M. Scanlon	18
		3.3.3	Aside: Philosophy's impact on Computing Science	19
3.4 Comparing Trust and Responsibility		Comp	aring Trust and Responsibility	20
	3.5 What work is missing?		work is missing?	20

4 Proposed Approach			Approach	21
4.1 A responsibility formalism's constituent e			onsibility formalism's constituent elements	21
		4.1.1	A trust formalism	21
4.2 Literature influence on the formalism's elements		ture influence on the formalism's elements	22	
		4.2.1	Formalism designed like Marsh's	22
		4.2.2	Outlining the proposed formalism	23
4.3 In answering research questions		wering research questions	24	
5	Wor	k Plan		25
5 Work Plan			20	

1 Introduction

Computational Responsibility is a field with little to no existing literature. Rather than a focus on *responsibility*, researchers have so far tackled a variety of other social topics through computational formalisation:

- Marsh's seminal work on Trust[Mar94]
- Stricter formal definitions on Trust, from a cognitive standpoint[CF]
- Some responsibility modelling, from a logical formalisation[SS15]

Finish reading this!

- Some work on reputation [CE11]
- Models of computational comfort models[MBEK⁺11].

While there is no direct literature on responsibility formalisms, then, we can see that there exists a wealth of literature for a responsibility formalism to be inspired by.

A responsibility formalism is useful in the same ways that formalisms of human traits such as reputation and trust might be; however, a responsibility formalism has the potential to have impacts in areas trust and reputation might not. For example, imbuing an intelligent agent with a sense of responsibility might provide it a greater degree of corrigibility[SFYA15]. An agent overseeing network security which understands its responsibilities within a much larger security system might better prioritise its duties when confronted with an unusual situation. Computational responsibility frameworks might help better model the emergent phenomena in sociotechnical systems, combine with traits like trust and comfort to make a more anthropomorphic device for better HCI, or perhaps help predict human actions in large computational models of human actors. We will explore some of these practical applications in § 4.

However, it is certain that a uses for these formalisms present themselves at every turn.

1.1 An early rebuttal of some common criticisms

One easy criticism to make of these anthropomorphic formalisms is the argument that, say, a trust formalism doesn't represent "true" trust. To address this point early, a responsibility formalism such as the one proposed need not be an entirely human-like representation of responsibility for every definition. Rather, there is a utility in an agent giving the *appearance* of responsibility. (If one follows the deterministic school of thought, there is also an argument that there is no difference[Hon].) Whether one considers it "true" responsibility should arguably be secondary to whether responsibility-like traits are useful to have computational frameworks for; we will see that these traits are indeed useful, and so that the criticism is moot.

Write something concrete here regarding fifelds it might be applicable in, like decision theory or AI safety or network security or socitechnical modelling

1.2 Terminology

Fill with terminology, in a style similar to honours dissertation

2 Statement of Problem

Computational responsibility is a complex area with lots of incidentally related work, but no specific relevant literature. Instead of focusing on the responsibilities of artificial agents, their responsibilities are implied by the construction of the agent itself. It might employ algorithms for driving without human guidance, or classify network traffic in an attempt to flag attempts at a system's security. In these instances, lots of somewhat-related work has been done on computational *trust*: can one artificial agent trust another?

However, simply building "responsibility" into a system without any understanding of responsibly-made decisions, or ability to learn obligations and duties in a concrete way, means we lose a great opportunity. Moreover, while trust and responsibility are intrinsically linked social concepts, no work has been done to migrate the models of trust to new models of responsibility that consider topics like obligation and duty. A concern arises: do artificially intelligent agents, which we put at the helm of concerns like network security and road safety, miss out as a result of their failure to consider duty and obligation? Two examples present themselves.

The first: a car might drive along a residential street and identify a squirrel running across the road in front of it. It calculates a high probability that, unless it swerves out of the way of the squirrel, it may kill it. It simultaneously identifies that, in the country it is driving in, the law states that it should swerve to avoid killing animals if possible. Computational responsibility introduces itself into the problem in that the car should also have a social understanding: will the swerve endanger humans? How strongly should it weight that probability into the action it chooses? Is it also responsible for, say, conserving fuel for environmental reasons? And if so, which responsibility is more important?

The key here is that the car has many goals to ascertain; while some are more immediate than others, it should have the capacity to weigh *multiple*, *arbitrary responsibilities* up to surmise what its next action is. Clearly, this is a problem for decision theory; but one where an understanding of responsibility may be of great help. Unfortunately, this example may be only expository as of yet: a practical model of a self-driving car with this degree of responsible awareness would be rather complex, and a model this advanced is beyond the scope of this project.

The second: two intelligent agents raise or lower the price of a book they manage according to percieved changes in the market. This is common practice on large sites where prices of unusual books can fluctuate according to sudden rises in demand, as seen in $\S 2$.

Confirm figure referencing style

In § 2, one artificial agent is known to have artificially inflated the price of a book; another agent has *also* inflated the price according to the seeming market trend. The first agent, seeing that the book is rising in value and now underpriced, inflates the price of its own copy, and the cycle continues until a human intervenes.

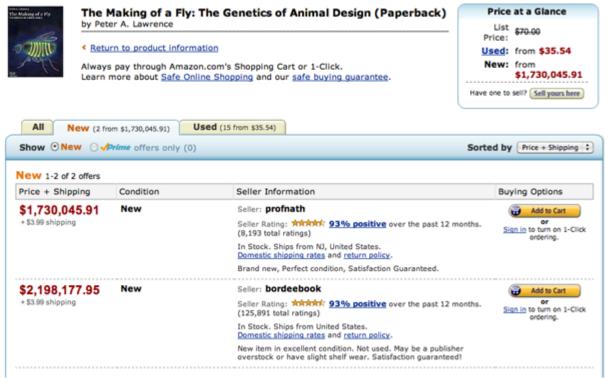


Figure 1: Bots on Amazon artificially inflate a book price to around 62850% its used price

Kevin Slavin discusses the idea that we have begun to design a world "for algorithms, with nothing but a big red button, labelled "stop" "[Sla]. The precession of this design trend marches on, relentless — but algorithms, rather than their interfaces, can be built with humans in mind. A mutual understanding of responsibility would allow one algorithm in this cycle to delegate the price inflation of its book to the other, breaking the cycle, so long as the concept of responsibility for a task is mutually understood. This is where the second, real-world example of computational responsibility lies.

As can be seen from the existence of models for concepts such as Trust which can solve similar HCI problems [PMB05], mimicking human traits computationally has its benefits. Moreover, we can be certain that, just as with trust modelling, useful and thorough responsibility models can produce work in machine ethics[Moo06], sociology[MW02], and clearly, computing science. We can therefore expect that a computational model of responsibility will yeild similar results — perhaps breaking new ground in other fields, which traits like Trust or Comfort have less relevance to.

2.1 Reflective Agents

As we move into a world increasingly dominated by algorithms and shaped by their decisions, there is a clear requirement for responsible systems. One problem arises: how can we be certain that an algorithm's internal conception of responsibility is 'human-like'? Early work by Sloman describes the notion of a 'space' of minds[Slo84], and this

concept is useful here. An artificial mind need not be human-like, or even biological-like; it can occupy an entirely different area of the space of minds altogether.

Therefore, when developing the proposed model of responsibility, one has to wonder what the components of the machine mind would be, such that it could house some useful definition of "responsibility". This useful definition need not be accurate; however, it will require the emulation of fundamental human attributes in order to successfully simulate.

For example: C&F define a "cognitive" agent as the lower limit of an agent's requirements for human traits for trust. C&F define a cognitive agent as:

Only a cognitive agent can "trust" another agent; only an agent *endowed* with goals and beliefs.

This definition doesn't quite fit our purposes — as will be seen, our definition also requires the concept of *obligation*. However, it can be seen that this definition is deliberately high-level in order to simulate the important components of a human trusting agent. A cognitive agent can be seen as an agent which, for the task it is set out to do, is modelled in a *high-level*, *human like way*.

Therefore, we might define our own high-level requirement of responsible computational agents:

Only a reflective agent can be "responsible" for its actions; only an agent which can *reflect on its obligations when choosing an action*.

This definition of a "responsible agent" as a "reflective" agent is important, because when considering its obligations, a responsible agent should be able to gauge whether to act in a certain way, weighted by their responsibility for a given obligation's fulfilment. As the model of responsibility developed begins to take shape, necessary components of those obligations — the responsibility equivalent of trust's goals and beliefs — will come to light.

2.2 Interpretive Agents

Unfortunately, simply reflecting on one's responsibility is not the only high-level requirement we can forsee needing for a responsible computational agent. Humans do not simply reflect on their obligations and duties before deciding on what their next actions might be. Human agents also see those obligations and duties through their own lens; they interpret their responsibilities according to certain factors which may influence their "feeling" of obligation.

Better examples?
Examples of both subjective interpretation and comparing two subjective interpretations? (Some agents might perceive one responsibility with the same score as more important than another responsibility, another agent might get it the other way around)

One can see this, for example, in people's respect for law or social convention. Sime citizens of a community might feel that it is imperative not to ride a bicycle on the pavement in Britain, as it is technically illegal. Others may well avoid the road traffic by making use of pedestrian areas if there aren't many pedestrians allowed, regardless of the law. Another example might be crossing the road; if a small child is present, parents of the child may well be teaching it to cross the road safely. To cross at a "red man" then, regardless of the presence of cars, somewhat derails the parent's lesson. It may even give the child an example of why they should be allowed to cross the road when they want. To cross the road at a "red man" does not respect one's influence over the situation at hand; in other words, the *responsible* thing to do is to wait for the lights to change.

However, it is clear that not everybody thinks this way; many cyclists ride on the pavements, and stopping at a red light to aid a parent in teaching their child might be considered extreme by some. The subjective nature of responsibility belies its interpretive nature: human actors interpret their obligations according to their beliefs, knowledge, and preferences, amongst other things. Therefore:

Involved in a "reflective agent" s judgement of their responsibilities is a subjective component: an interpretive function which converts information about an obligation or duty into a subjective score of responsibility.

This way, the human-like subjectivity of responsibility can be simulated.

2.3 Reactive Agents

There is one last definition important to our understanding of what a basic responsibility formalism might look like. If we want to represent a human-like sensation of trust, it's important that we acknowledge that human agents become more or less trustworthy over time. Children, for example, can be trusted much more once they mature and become adults. Moreover, that adult agent will probably be slightly less likely to trust than it was as a naive child. Therefore, to simulate human-like responsibility, agents' interpretation of responsibility should change over time.

As we'll see, responsible agents whose sense of responsibility changes as a reaction to its environment is a notion which is backed not only by observation, but by the moral philosophy which underpins the proposed work at the end of this review. For now, we can define these "reactive agents" as agents whose sense of responsibility shifts as their environment is seen to change:

Only a "reactive agent" has a changing subjective outlook on the world; it changes its reflection on its own and other agents' responsibilities depending on its environment.

It is clear to see that for a reactive agent to change its reflection on responsibility, it must be an interpretive agent; in other words, the set of reactive agents is a subset of the set of interpretive agents. It should also be clear to see that a useful responsible agent should be both *reactive* and *reflective*: the set of reactive, reflective agents have a perception of the world which is subjective, can change according to its environment, and uses its sense of responsibility when making decisions.

2.4 Using the definitions

These definitions are useful in that they allow us to begin to see what might compose a responsible agent. We might begin to build a computational model of responsibility which embodies these traits by investigating the following research questions:

- 1. How can a computational formalism of responsibility direct the decisions made by an intelligent agent?
- 2. How can an intelligent agent assume the consequences of actions it makes, the decisions other agents make, and its general environment, so as to direct its interpretation of responsibility?

These are the research questions I seek to answer in the course of this project.

3 Background Survey

Unlike Computational Responsibility, Computational Trust is a topic which has a surprising degree of pre-existing literature. Marsh [Mar94] draws inspiration from as early as David Birkhoff's 1930s work in creating an 'Aesthetic Measure', where the famous mathematician created a quantification of Aesthetics. While some dispute that such subjective topics can be boiled down to a single number (or array thereof), much work to the contrary has now been completed. Like Marsh, we should start from the beginning.

3.1 Social Sciences and Mathematics

3.1.1 Birkhoff's Aesthetic Measure

One of the earlier formalisms of a human factor¹² was Birkhoff's definition of Aesthetic Measure[Bir]. In it, Birkhoff defines the notion of Aesthetic Measure as a ratio of Order to Complexity:

$$M = \frac{O}{C}$$

Birkhoff's work inadvertently gave rise to the notion that human factors can be represented by mathematical equations and systems. Birkhoff's formalism of aesthetics became popular for a few reasons, but one of particular interest to later Trust modelling work was that Birkhoff put a great degree of effort into backing his work up with psychological theory. In this way, Birkhoff's formalism could be said to be a *psychological* formalism.

Later trust modelling work followed in Birkhoff's footsteps here. Indeed, Birkhoff gives a solid foundation for the model-creating method later employed by Marsh[Mar94] and Castelfranchi & Falcone, as it is:

- Founded on mathematical or logical principles which are quantifiable
- Heavily inspired and directed by related work in psychology, sociology, and philosophy

- Humans take orders and manage the running of the shop
- Technology is responsible for complex activities such as taking payments and forcing steam through coffee at high pressure

so there are both social and technical actors and behaviours in the "system" of a day-to-day coffee shop.

¹For the sake of clarification, we define a "human factor" as an element of a social or sociotechnical system which arises from human behaviour, such as Trust.

²Also for the sake of clarifying a sociotechnical system, a sociotechnical system is a system composed of human tendencies and behaviours, such as Trust, alongside technical activity, such as a computer or a steam engine. An example might be a coffee shop:

CITE THIS

The marriage of social studies with mathematical rigour will be a recurring theme of the work related to Computational Trust.

3.1.2 Deutsch

Following the quantifiable, mathematical work done by Birkhoff, logical and arithmetic formalisms of human factors followed. One of the earlier and more widely adopted models for Trust came from Deutsch in 1962. Deutsch is a psychologist who did swathes of work in the topic of cooperation, touching on Trust during the 60s.

Deutsch's formalism of trust wasn't immediately quantifiable, but presented one of the earliest well-defined definitions of trust. To paraphrase Deutsch's formalism in "Cooperation and Trust: Some Theoretical Notes", 1962:

- An actor is presented with a choice between two paths.
 - A: No change
 - *B*: The actor takes some action, of ambiguous outcome. A possible gain is associated, *P*, and some possible risk is associated, *R*.
- The actor assesses that the outcome of choice *B* relies on the behaviour of another actor.
- The actor assesses the action they may take and resolves that the strength of *R*, likelihood of *R* as an outcome, or both are higher than the respective *P* measurements.
- The actor is said to be *trusting* they take path *B*.

This formalism introduces some interesting notions. For example, it is unclear as to whether the outcome of choice *B* can rely on the same actor making the decision; can one trust oneself by Deutsch's definition? Another interesting analysis of the implications of Deutsch's model is that it does not rely on the *accurate* measurement of risk and utility, but just its perception — trust is subjective, and based on the trusting actor's perspective on the world.

Rather than characterising trust by the parties involved, Deutsch's formalism is characterised by *risk and utility*. A simple quantification of Deutsch's formalism could be devised, therefore, where risk and utility are quantified by simple assessments using utility functions and a form of risk analysis. Even so, the outcome of this quantified system is a single bit: trusting or not trusting. This does quantify trust, but only technically speaking, and this quantification is weak in its expressiveness. It gives no remit to suggest that one might trust one person over another, for example, as there are no orderable degrees of trust.

Deutsch offers many different ideas as to why and how trust or trust-like behaviour can come about, however. This list is taken from Marsh 1994[Mar94], where explanations of all nine can be found:

- 1. Trust as Despair
- 2. Trust as Social Confirmity
- 3. Trust as Innocence
- 4. Trust as Impulsiveness
- 5. Trust as Virtue
- 6. Trust as Masochism
- 7. Trust as Faith
- 8. Risk-taking or Gambling
- 9. Trust as Confidence

Deutsch's given model above specifically targets formalisation of trust as confidence.

3.1.3 Luhmann

Luhmann, a sociologist who also worked in Trust and related fields, had his own take on formalisms of Trust: that trust was a social tool for reducing the complexity of a social system. Specifically, Lohmann sees trust as being a method whereby agents in a social system can reduce their exposure of *risk* to each other. According to Luhmann, "Trust… pressuposes a situation of risk."

CITE

Luhmann's work is therefore difficult to form quantitative formalisms from, as his thesis stems from a risk analysis perspective, which can be particularly difficult in a sociotechnical system. However, Luhmann's work remains interesting; a formalism of a human factor like trust would be incomplete without considering the properties of individual human actors as well as these properties' emergent effects in the larger sociotechnical space. For small systems, these social-level properties may not present themselves very strongly; however, most human factors are present regardless of the scale of the system being modelled. Therefore, a formalism of a human factor which fails to consider both psychological and sociological aspects cannot be complete.

3.2 Modern [Computational] Trust methods

3.2.1 Marsh's formalism

The earliest quantifiable formalism of trust which provides computability, flexibility, and an inspiration from the sociological and psychological work above is that of Stephen Marsh in 1994[Mar94]. Marsh's work breaks trust up into three core quantifications, where each variable takes some value in the range [-1,1):

1. Basic Trust

This is the general degree of "trustingness" about an agent, or that agent's ordinary inclination to trust.

2. General Trust

General trust is trust in the context of the agent being trusted. Marsh's original description begins[Mar94]:

Given two agents, $x, y \in A$, to notate 'x trusts y' we use: $T_x(y)$ The value represents the amount of trust x has in y here.

So, General Trust can be seen to be the trust that an agent *x* has in *y*.

3. Situational Trust

Trust doesn't exist in a vaccum, and the only variable isn't the subject of x's trust; y may have varying degrees of competency in performing an action. Therefore, Situational Trust can be seen to be the trust x holds that y can actually perform some task, α . Marsh helpfully gives the example[Mar94]:

... whilst I may trust my brother to drive me to the airport, I certainly would not trust him to fly the plane!

Marsh's three types of trust are helpful in breaking down what matters when discussing trust — notions like competency, for example — as well as establishing a jargon for trust. Often, one might say that a person is "trusting": Marsh's formalism accounts for concepts like this, but establishes it as a less detailled type of trust, and a type of trust which doesn't account for the action being trusted for, or whether the trusted agent is able to complete the action.

Marsh also succeeds in introducing concrete examples of computational formalisms of ordinarily human traits — here Trust. The key aspect of Marsh's advancement is that it goes one step further than a *quantitative* model, and introduces reinforcement learning algorithms which model how trust *changes*, and not just its current state. As seen when discussing Birkoff's work (§ 3.1.1), quantitative formalisms of human traits like Aesthetics had been studied and achieved long before Marsh's work.

Since Marsh's work, many trust models have been developed. A small subset of these are reviewed here; offshoots from this seminal work include REGRET, Eigentrust, FIRE,



and others.

Finish notes on Marsh's computational trust model!

3.2.2 Castelfranchi & Falcone

As it turns out, cognitive computational trust models that already exist are almost but not quite appropriate for modelling responsibility. The C&F trust model requires only four main ingredients to formulate a cognitive trust model:

- 1. x, a truster
- 2. *y*, a subject of trust
- 3. *g*, a goal of *x*
- 4. α , an action of y

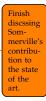
This model gets us close to where we need to be to model responsibility; like responsibility modelling often does, it assumes two agents. There also exists some goal which can be met, which — to use C&F terminology — is *delegated* by x to y. Y can achieve this goal through some action, α . So far, all of this forms the beginning of a foundation for cognitive responsibility; what turns delegation of a task into the consignment of responsibility is that of obligation, and the understanding of obligation.

It is evident that trust and responsibility models are, even in the human-like cognitive approach, very similar. However, there are drawbacks which mean that we cannot directly apply C&F theory to the idea of computational responsibility: it does not represent any degree of obligation or address the specific problem of judging responsibility at all.

Nevertheless, this presents an exciting insight into work to be done to produce a formalism of responsibility. Particularly, it is evident that there is at least some technical value in listing the individual components as C&F do. Their simple, reduced approach implies that with the correct identification of elements of responsibility, our formalism can be similarly simple. It is also enouraging that connections between trust and responsibility modelling seem to readily present themselves. We can therefore expect our formalism to rightly exhibit a similar structure and features.

3.2.3 Ian Sommerville, Sociotechnical Systems, and Responsibility Modelling

Sommerville's work focuses largely on sociotechnical systems and responsibility modelling — in this way, Sommerville's work is not typically concerned with computational models of trust, as the above were. However, his work does begin to border on our own advancements, providing responsibility modelling formalisms.



It is important to note that Ian Sommerville has been a particularly prolific writer for a researcher in the sociotechnical systems scene. Sommerville's modelling systems are sometimes graphical[LSSB10]. Unfortunately, graphical modelling systems do not lend themselves particularly well to computational formalism: they don't yield naturally to numerical analysis; they are generally designed for the purposes of human visual analysis, instead of logical reasoning; they are often difficult to represent non-graphically, which arguably makes input and manipulation too complex for the purposes of designing a complex intelligent agent around.

Ultimately, though, graphical responsibility modelling systems are designed for representing the responsibilities of a single agent at a given point in time; a responsibility formalism, by contrast, should be a series of metrics and rules which can apply to arbitrary reasoning of an agent's responsibility through time, with that agent using the formalism to reason about its changing responsibilities. In other words, graphical responsibility modelling differs from a responsibility formalism in that a responsibility formalism needs to *generalise* reasonsing about how responsibility, as a concept, "behaves".

Nevertheless, some sociotechnical work on responsibility prevents an invaluable addition to relevant literature for developing its computational formalism. In particular is Sommervile's work on "Causal" and "Consequential" responsibilities. In defining these terms, Sommerville writes[Som07]:

Consequential responsibility can only be assigned to a person, a role or an organisation automated components cannot be blamed. Causal responsibility reflects who or what is responsible for making something happen or avoiding some undesirable system state. It is often the case that these are separated.

The seperation of concerns between consequential and causal responsibilities can help us to inform the structure and nature of a responsibility formalism. Particularly, one can simplify these definitions to be that:

Consequential responsibilities are responsibilities which relate to a state arrived at in the past, and the relationship of actors and actions with said state.

Causal responsibilities are responsibilities which relate to future states, and the actors and actions which are potentially assigned with the goal of arriving at said state.

This seperation of past repsonsibilities and future responsibilities means that we can structure the concerns and operation of a responsible agent according to a given formalism. Particularly, for our purposes, the assessment of one's responsibility to arrive

at future states — one's assessment of its *Causal Responsibilities* — might be informed by information an agent's information about its previous responsibilities and its ability to act responsibly — one's *Consequential Responsibilities*. Therefore, we might decide to create a responsibility formalism which specifically models the change in causal responsibility, by assessing consequential responsibility.

As we approach a review of philosophical literature in § 3.3, we will find that a formalism geared toward causal responsibility is not only a very attractive way to model from a philosophical perspective, but that the philosophical literature on responsibility has converged on similar definitions of responsibility. Therefore, there appears to be a strong case that causal-first models — aside from their simple structure and readiness to be tied into an agent's decision functions — are a sound way of approaching the problem of framing problems concerning responsibility, because of the consensus across fields which rarely interact.

3.3 Philosophy of Moral Responsibility

Philosophy regarding moral responsibility is an area whose literature is both wide and deep. That said, not all moral repsonsibility literature is relevant to a computational repsonsibility project; lots of it is designed from a social analysis perspective which would be difficult to implement in any useful way. Other areas, however, present more promise for studies regarding formalisms.

3.3.1 Peter F. Strawson

One example of research with utility in a computational way is that of Peter F Strawson, particularly in his seminal essay, Freedom and Resentment [Str62]. Strawson's topic actually revolves around whether determinism has any impact on "free will" — a discussion clearly outside the scope of this project — but in forming his argument creates some key conepts that we can use to consider the applicability of a responsibility formalism to a computational system, as well as touching on what that formalism would look like.

Strawson's fundamental argument can be construed as being that determinism doesn't affect what human factors — like Responsibility and Trust — mean, because these concepts are founded on the relationships between human actors, rather than being inherant to the human actors themselves. As computer scientists, we can extrapolate this argument out to a sociotechnical environment. That is: Strawson's argument applies to both socila *and* technological agents within a sociotechnical world. Using Strawson's reasoning, then, we can firmly conclude that it *does* make sense to create "responsible" computers, because responsibility makes sense as a trait for a computer to have in the same way it might a human actor.

Strawson's insights don't stop there, though. He also produces an interestingly rigorous analysis of how ordinarily fuzzy human factors can be formalised appropriately:

Indignation, disapprobation, like resentment, tend to inhibit or at least to limit our goodwill towards the object of these attitudes, tend to promote an at least partial and temporary withdrawal of goodwill; they do so in proportion as they are strong; and their strength is in general proportioned to what is felt to be the magnitude of the injury and to the degree to which the agents will is identified with, or indifferent to, it.

[Str62]

Strawson captures an essential component of Marsh's computational trust model: an actor's actions influence the percieved trustworthiness of an agent in proportion of the magnitude of the "trustworthiness" of their actions. This concept, Strawson's elucidation shows, is one which can also extend to responsibility; it influences all of the human factors Strawson seeks to address in his essay (i.e. all human factors).

This gives us an insight into how a realistic computational trust model might function: its perception of responsibility should alter depending on external factors. We are aware, however, that people's actions are not the only things which can affect our feeling of responsibility: all sociotechnical factors might. For example, the "Bystander Effect" [FGPF06] occurs when one feels less responsible for acting in an emergency situation because of the number of people present who could help; in the end, nobody does, because every person's sense of responsibility is weakened. However, it is not weakened by any person's actions: it's influenced by a set of sociotechnical factors, of which the feelings that Strawson discusses are a subset.

A suitable conclusion one might draw, then, would be that a valid computational model of responsibility *must* take into account an analysis of the sociotechnical environment of the responsible agent. It would also be valid to draw the conclusion — as a result of the first of Strawson's points discussed — that it makes sense to talk about "trusting" and "responsible" computational agents as if they were humans. Another important conclusion to draw from this part of Stawson's discussion is that this is one way in which agents' interaction is meaningful: therefore, there is a significant body of work to be undertaken in the combination of trust and responsibility formalisms to the field of HCI. Some of this work is already underway[MBEK+11, MW02].

Strawson also notes what sort of agent one rightly considers responsible. His argument relates to which agents one judges responsibility *in*, rather than judging which actions an agent might consider itself responsible for. However, it is pertinent to this research in that it highlights what assumptions the formalism might make about an agent it models the responsibility of. In other words, Strawson helps us to limit the scope of the formalism in terms of the agents it targets.

Let us consider, then, occasions for resentment ... To the first group belong [agents of whom we can say] "He didn't mean to", "He hadn't realised", "He didn't know" ... "He couldn't help it" ... None of them invites us to suspend towards the agent, either at the time of his action or in general, our ordinary reactive attitudes.

The second and more important subgroup of cases allows that the circumstances were normal, but presents the agent as psychologically abhormalor as morally undeveloped. The agent was himself; but he is warped or deranged, neurotic or just a child. When we see someone in such a light as this, all our reactive attitudes tend to be profoundly modified.

[Str62]

Strawson here demonstrates a useful seperation of two groups of agents: those who incur resentment through a lack of control but who are fully able to act correctly, and those whose lack of control or basic understanding disqualifies them as agents who can incur resentment at all. A definition of resentment would be useful here; unhelpfully, Strawson doesn't lend one, but helpfully, a reasonable general definition of resent is easy to formulate: one feels resentment toward another agent when a goal entrusted to it is not achieved.

In other words, one feels resentment toward an agent which *fails to fulfil a responsibility*. Using this definition, we can use Strawson's separation of resentment-inducing agents to limit the responsibility formalism's scope: an agent can be considered responsible for their actions when they can reasonably be assigned some goal (and possibly some actions to achieve this goal). Another way of saying this would be that an agent can be considered responsible for an action if we can *trust* the agent with a goal; if a goal is assigned, then the act of assignment makes the agent responsible.

Using Strawson's reasoning regarding resentment, then, we can reasonably assert the implication: a responsible agent is an agent actively trusted with a task: a "causal responsibility", to borrow Sommerville's terminology. Along with the earlier notes on types of agents § 2.2, we can even further limit the scope of agents the formalism should be targeted toward. We also neatly tie into our formalism the assignment of responsibility to an agent; more detail, along with what the assigned responsibility should represent, is discussed in the proposed approach § 4.

3.3.2 Thomas M. Scanlon

In his essay, Justice, Responsibility, and the Demands of Equality[Sca06], Thomas Scanlon assesses responsibility as having two factors which bear striking resemblance to Sommerville's "Causal" and "Consequential" Responsibilities: "Attributive" and "Substantive" Responsibilities.

Scanlon discusses "Attributive" responsibilities as a group of duties according to the definition:

What a person sees as a reason for acting, thinking, or feeling a certain way.

For the purposes of later discussion, utility, and slight simplification, I will generalise Scanlon's definition to be "responsibilities for future actions". He also discusses "Substantive" responsibilities, which are loosely defined as responsibilities for the choices an agent has made, taking into account the effects they have and obligations at the time. We can generalise these to be "responsibilities for past actions". While in doing so, I reduce Scanlon's definitions to a discussion of action as opposed to the wider gamut of human properties, this framing sufficiently limits the scope of the work to apply elegantly to decision theory. Therefore, we can begin to apply Scanlon's work to the actions taken by responsible computational agents.

This limitation of scope also enables us to tie Sommerville's sociotechnical research to work done in the philosophical sphere. A complete responsible computational system, therefore, would be suitable for the same scrutiny that human agents currently undergo in the field of moral philosophy. It also, like Sommerville's work, presents us with a framework for thinking about the scope of responsibilities that an artificially intelligent agent might be subject to.

Beyond philosophical waxing lyrical, we have reason to limit our formalism to a subset of responsibilities Sommerville and Scanlon describe. Specifically, we can use Sommerivlle's "Consequential" or Scanlon's "Substantive" responsibilities to tailor a responsible agent's judgement of its "Causal" or "Attributive" responsibilities. This structure will act as a scaffolding for the simplified proposed framework in § 5.

3.3.3 Aside: Philosophy's impact on Computing Science

It's worth noting that Philosophy and Computing Science are unusual bedfellows. However, this need not be the case. As we've seen already, and will continue to observe, there is a definite need for computing science to make use of and contribute to research involving human factors — whether those factors be sociological, psychological, ethnographic, or even philosophical. What's more, philosophical research may well have an impact on the ethics and cultural implications that certain advances in Computing Science might entail. When discussing matters of potentially existential levels of risk, it is important to have the foresight to construct solid, informed arguments about our role in the grander scope of things.

To this end, I believe collaboration between Philosophy and Computing Science should be a more routine affair. Some interdisciplinary work of this form is already underway: Nick Bostrom's Future of Humanity Institute at Oxford University and the Machine Intelligence Research Institute in the USA are two excellent examples of Computing Science and Philosophy coming together to produce vitally important work, which greatly influences the progression of the field of AI Safety.

3.4 Comparing Trust and Responsibility

Trust and responsibility are similar concepts. Seen through the lens of C&F, this is particularly clear:

- They both concern themselves with actions and goals
- They're both naturally framed in terms of task delegation
- They both follow Strawson's note on how agents with human factors change their outlook with respect to these factors through interaction.

Trust and Responsibility's similarity regarding task delegation is particularly interesting: this might allow a responsibility formalism to operate in a similar way to trust, meaning that much of the existing literature on Trust could be reappropriated (with some research) for the purposes of computational responsibility. Indeed, one definition of responsibility for a task might be the obligation to act on a delegated task. A corollary of C&F argue that task delegaton is inherantly trust; then, all one need do to extend a trust formalism to a responsibility formalism would be to augment C&F's axioms to include an agent's obligation, and one's formalism would be complete!

Unfortunately, C&F's formalism, while technically calculable by a computer, uses logical expressions to evaluate trust. For the purposes of an application such as a decision function of an intelligent agent, unless that agent follows a structure such as a BDI model, this would not be terribly useful — though it's worth noting that granular models of C&F also exist[].

find citation for this

3.5 What work is missing?

Develop arguments that the responsibility formalism might actually be put to good use, as per 8.1

Discuss responsibility assignment through

Discuss the interpretation of responsibility through the interpretation function — relate to Marsh's specific trust, and how it's interpreted?

4 Proposed Approach

We can see that the literature available, while not on computational responsibility formalisms, are all fairly related to constructing the formalism proposed. In light of this, we can begin to identify some of the components of a suitable formalism:

- The formalism should answer the research questions laid out in the problem statement § 2:
 - 1. How can a computational formalism of responsibility direct the decisions made by an intelligent agent?
 - 2. How can an intelligent agent assume the consequences of actions it makes, the decisions other agents make, and its general environment, so as to direct its interpretation of responsibility?
- The formalism should suitably limit the scope of the actors it models to be reactive, reflective agents § 2.2.

 The formalism doesn't apply to all agents; we can imagine some agents which are *not* responsible, such as those shown by Strawson § 3.3.1, and agents which aren't reactive or reflective.
- The formalism should interpret the obligation an agent is assigned when another agent trusts it with a causal responsibility.
- Ideally, the formalism would utilise as much relevant psychology and sociology literature as possible so as to maximise the formalism's interdisciplinary potential.

4.1 A responsibility formalism's constituent elements

4.1.1 A trust formalism

A useful exercise is to discern what a responsibility formalism would consist of. Worth noting is that this formalism is a proof-of-concept and a starting point for further research to refine. Therefore, a simple model which is easy to implement correctly and reason about should be paramount.

With this in mind, one important constituent part is that of the trust model the formalism works in tandem with. While this could in theory be any model, three notions are worth bearing in mind when selecting one.

- 1. The trust formalism selected should act as a starting point for the responsibility formalism's design, because of the similarities between trust and responsibility.
- 2. The trust formalism should express gradations of trust, rather than a boolean logical approach. This is important as the trust an agent has in another agent will allow the former to assess how responsible the latter is; these agents may be required to make judgements of *how* responsible other agents may be, so they might choose one agent to be responsible for a goal over another.
- 3. The trust formalism should be simple, to act as a basis for the design of a simple responsibility formalism.

To satisfy all of these aims, Marsh's seminal trust formalism seems most appropriate. This is for a number of reasons. For example, Marsh's formalism is very easy to understand and implement — particularly with it being an early formalism uncomplicated by later literature. Another reason is that Marsh's formalism allows one to express gradations of trust; in contrast to another model, such as Castelfranchi & Falcone's model, Marsh's formalism requires no additional complexity to model trust gradation.

Marsh's model is also constructed with other disciplines in mind — psychology and sociology feature prominently in its cited literature. However, Marsh cites no philosophical advancements in his model; therefore, the application of moral responsibility to the formalism being designed cannot be based on similar work used with this formalism.

4.2 Literature influence on the formalism's elements

4.2.1 Formalism designed like Marsh's

The model of responsibility being designed requires the ability to model gradations of trust; therefore, Marsh's Trust formalism will act as the template for the responsibility formalism's design. In particular, Marsh's model's segregation of trust are useful for establishing the basic principles behind the responsibility model:

Marsh's Trust	Responsibility Formalism
Basic Trust	Basic Responsibility
General Trust	General Responsibility
Specific Trust	Specific Responsibility

Using Marsh's model and adopting the jargon he develops means we can keep lots of the notions he develops, such as variable domains and separation of different "levels" of responsibility, such as how generally responsible an agent is, and how responsible an agent is for a very specific thing. Indeed, a layout of the variables we might use in a formalism of responsibility, compared to Marsh's, might look like this:

Marsh's Trust Model Variables		
Knowledge (of x at time t)	True/False	
Importance (of knowing fact x at time t)	[0, +1]	ם
Utility (of action α at time t to an agent)	[-1, +1]	General
Basic Trust (of agent A at time t)	[-1, +1)	
General Trust (of agent A at time t in agent B)	[-1, +1)	Situational Resp
Situational Trust (of agent A at time t in agent B doing action α)	[-1, +1)	

Marsh's formalism also makes use of other concepts, such as sets of agents and definitions of what trust might be composed of; similarly, responsibility modelling will require the development of more concrete definitions, as we will see.

Outlining the proposed formalism

An example breakdown of responsibility and what it is might be the following:

Responsibility Term	Definition of Term
Agent / Actor	Combination of Constraints, Environment, and Be-
_	liefs. Acts on a decision function.
Obligation	Combination of Authority, Goal, Set of Appropriate
	Actions, Responsibility Score
Responsibility Score	Number in (0,1] assigned by authority denoting de-
	gree of obligation
Action	Combination of Activity, Resource requirements, Pos-
	sible Issues and Effect
Authority	Agent which assigns an obligation to another agent

Not all terms above are completely defined; however, the table is provided as an example of how a responsibility formalism might look. The formalism proposed here is by no means final, but it can be seen that it would produce gradations of responsibility, that authorities assign an obligation with a certain score, and that agents choose actions based on a decision function — a decision function which takes into account assigned obligations when it produces a next action.

One difference between the above table and a final formalism of responsibility is that the final formalism would also account for processes. For example, a number assigned by an Authority represents the degree to which the responsible agent is obliged to act toward a goal, but the agent itself needs to weigh this up against other factors. For example: might the goal be time sensitive? It's imperative that I pay my rent on time, but a responsibility to keep windows closed doesn't depend on any one point in time. In other words, regardless of the initial score assigned to paying rent, I (as an agent) must interpret that score while taking into account other features of the goal — these features change from goal to goal, and theoretically from agent to agent, too. This interpretation function is a necessary part of the formalism, too — but it is defined by process, and not semantically.

4.3 In answering research questions

In stating the problem proposed, two research questions were devised:

- 1. How can a computational formalism of responsibility direct the decisions made by an intelligent agent?
- 2. How can an intelligent agent assume the consequences of actions it makes, the decisions other agents make, and its general environment, so as to direct its interpretation of responsibility?

The first research question can be answered in part by the definition of the Agent / Actor's decision function making use of the responsibility formalism in choosing the agent's next actions. To show this, artificial agents will be constructed with decision functions which are guided by the agent's assigned responsibilities.

The second research question can be answered with similar methods. To construct artificial agents with responsibility-guided decision functions, a full responsibility formalism will be devised and implemented in these artificial agents. In doing this, the formalism will work in tandem with the interpretation functions of the agents developed. In doing so, the interpretation function developed will show that a complete implementation of the responsibility formalism will answer the second research question.

It is worth noting that the interpretation function which answers the second research question is a necessary component of the formalism, but no one interpretation function belongs in the responsibility formalism itself. This is because, while a function can be imagined which has certain properties, two different agents may require different interpretation functions. In other words, the interpretation is a separate concern from the formalism itself: the formalism might specify *domains* for the interpretation function to map from and to, but cannot specify specific processes and parameters. Therefore, no canonical version of the interpretation function may be provided, though some properties of the function will be noted and incorporated into the formalism's definition.

5 Work Plan

Panic, write the report in a 36 hour caffeine-induced fever dream

Todo Notes

Finish r	eading this!	3
	omething concrete here regarding fields it might be applicable in, like on theory or AI safety or network security or socitechnical modelling .	3
Fill with	n terminology, in a style similar to honours dissertation	4
Confirm	n figure referencing style	5
two s bility	examples? Examples of both subjective interpretation and comparing ubjective interpretations? (Some agents might perceive one responsiwith the same score as more important than another responsibility, anagent might get it the other way around)	8
Clunky	? Re-word?	8
CITE TH	HIS	11
CITE TH	HIS	12
Should	we discuss this too?	13
Finish n	notes on Marsh's computational trust model!	14
Finish d	liscssing Sommerville's contribution to the state of the art	14
find cita	ation for this	20
	o arguments that the responsibility formalism might actually be put to use, as per § 1	21
Discuss	responsibility assignment through trust	21
	the interpretation of responsibility through the interpretation function ate to Marsh's specific trust, and how it's interpreted?	21
Referer	nces	
[Bir]	G D Birkhoff. AESTHETIC MEASURE.	
[CE11]	Partheeban Chandrasekaran and Babak Esfandiari. A model for a testh for evaluating reputation systems. <i>Proc.</i> 10th IEEE Int. Conf. on Trust, Serity and Privacy in Computing and Communications, TrustCom 2011, 8th IE Int. Conf. on Embedded Software and Systems, ICESS 2011, 6th Int. Conf. FCST 2011, pages 296–303, 2011.	еси- ЕЕЕ

- [CF] Cristiano Castelfranchi and Rino Falcone. Social Trust: A Cognitive Approach.
- [FGPF06] Peter Fischer, Tobias Greitemeyer, Fabian Pollozek, and Dieter Frey. The unresponsive bystander: are bystanders more responsive in dangerous emergencies? *European Journal of Social Psychology*, 36(2):267–278, 2006.
- [Hon] Ted Honderich. Free will, determinism and moral responsibility the whole thing in brief.
- [LSSB10] Russell Lock, Tim Storer, Ian Sommerville, and Gordon Baxter. Responsibility modelling for risk analysis. 2010.
- [Mar94] Stephen Paul Marsh. Formalising Trust as a Computational Concept. *Computing*, Doctor of(April):184, 1994.
- [MBEK+11] Stephen Marsh, Pamela Briggs, Khalil El-Khatib, Babak Esfandiari, and John A. Stewart. Defining and Investigating Device Comfort. *Journal of Information Processing*, 19(7):231–252, 2011.
- [Moo06] James H. Moor. The nature, importance, and difficulty of machine ethics. *IEEE Intelligent Systems*, 21(4):18–21, 2006.
- [MW02] Michael W. Macy and Robert Willer. From Factors to Factors: Computational Sociology and Agent-Based Modeling. *Annual Review of Sociology*, 28(1):143–166, 2002.
- [PMB05] Andrew Patrick, Stephen Marsh, and P Briggs. Designing systems that people will trust. *Security and Usability: Designing Secure Systems That People Can Use*, pages 75–100, 2005.
- [Sca06] Thomas M Scanlon. Justice, responsibility, and the demands of equality. 2006.
- [SFYA15] Nate Soares, Benja Fallenstein, Eliezer Yudkowsky, and Stuart Armstrong. Corrigibility. *AAAI Workshop on AI and Ethics*, (2014):74–82, 2015.
- [Sla] Kevin Slavin. HOW ALGORITHMS SHAPE OUR WORLD.
- [Slo84] Aaron Sloman. The Structure of the Space of Possible Minds. pages 35–42, 1984.
- [Som07] Ian Sommerville. Models for responsibility assignment. In *Responsibility* and dependable systems, pages 165–186. Springer, 2007.
- [SS15] Robbie Simpson and Tim Storer. Formalising responsibility modelling for automatic analysis. In *Lecture Notes in Business Information Processing*, 2015.
- [Str62] Peter F. Strawson. Freedom and resentment. *Proceedings of the British Academy*, 48:1–25, 1962.