



University of Glasgow | School of  
Computing Science

# On Computational Responsibility

William T. Wallis

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

Masters project proposal

# Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Introduction</b>   | <b>3</b> |
| 1.1      | An Early Rebuttal of some Common Criticisms . . . . .                   | 3        |
| 1.2      | The Scope of the Model . . . . .  | 4        |
| 1.2.1    | Reflective Agents . . . . .   | 4        |
| 1.2.2    | Interpretive, Reactive Agents . . . . .                                 | 5        |
| 1.3      | Proposal Overview . . . . .   | 6        |
| <b>2</b> | <b>Statement of Problem</b>   | <b>7</b> |
| 2.1      | Research by Scientific Method . . . . .                                 | 7        |
| 2.2      | Appropriately Fit Agent Context . . . . .                               | 8        |
| <b>3</b> | <b>Background Survey</b>  | <b>9</b> |
| 3.1      | Social Sciences and Mathematics . . . . .                               | 10       |
| 3.1.1    | Birkhoff's Aesthetic Measure . . . . .                                  | 10       |
| 3.1.2    | Deutsch . . . . .   | 10       |
| 3.1.3    | Luhmann . . . . .   | 11       |
| 3.2      | Modern [Computational] Trust methods . . . . .                          | 12       |
| 3.2.1    | Marsh's formalism . . . . .   | 12       |
| 3.2.2    | Castelfranchi & Falcone . . . . .                                       | 13       |
| 3.2.3    | Eigentrust . . . . .  | 14       |
| 3.3      | Ian Sommerville, Sociotechnical Systems, and Responsibility Modelling . | 15       |
| 3.4      | Philosophy . . . . .  | 17       |
| 3.4.1    | Peter F. Strawson . . . . .   | 17       |
| 3.4.2    | Thomas M. Scanlon . . . . .   | 19       |
| 3.4.3    | Deontic Logic . . . . .   | 20       |
| 3.4.4    | Sloman . . . . .  | 21       |

|          |   |           |
|----------|---|-----------|
| 3.5      | Discussion . . . . .  | 21        |
| 3.5.1    | The Relationship Between Trust and Responsibility . . . . . | 22        |
| <b>4</b> | <b>Proposed Approach</b>                                    | <b>23</b> |
| 4.1      | A Responsibility Formalism's Constituent Elements . . . . . | 23        |
| 4.1.1    | A Trust Formalism . . . . .                                 | 23        |
| 4.2      | Literature Influence on the Formalism's Elements . . . . .  | 24        |
| 4.2.1    | Formalism Designed Like Marsh's . . . . .                   | 24        |
| 4.2.2    | Outlining the Proposed Formalism . . . . .                  | 25        |
| 4.3      | In Answering Research Questions . . . . .                   | 26        |
| <b>5</b> | <b>Work Plan</b>  | <b>27</b> |
| 5.1      | Completing the Formalism . . . . .                          | 28        |
| 5.2      | Designing the Use Case . . . . .                            | 28        |
| 5.3      | Testing the Formalism . . . . .                             | 28        |
| 5.4      | Write-up . . . . .  | 29        |
| 5.5      | Review of Time Allocation . . . . .                         | 29        |

# 1 Introduction

Computational formalisms of social constructs are an increasingly common research area. For example, researchers have so far tackled a variety of social notions through computational formalism:

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

These models of social constructs are useful in a variety of ways: Marsh’s model, for example, gave rise to new methods in solving problems in fields as diverse as HCI[PMB05] and systems modelling[HJS06]. However, responsibility as a social construct has been neglected; no literature on responsibility formalisms has been published to date. This is curious, as responsibility modelling is a field which has proven particularly useful — therefore, a logical next step for responsibility as a subject of study within sociotechnical systems analysis would be the computational formalism of the trait. As will be demonstrated, responsibility as a computational concept may yield a great number of research opportunities in fields such as sociotechnical systems modelling, machine learning and decision theory, and even humanities such as the philosophy of mind.

A responsibility formalism is useful in the same ways that formalisms of human traits such as reputation and trust might be; however, a computational theory of responsibility has the potential to impact areas which 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; they might combine with traits like trust and comfort to make a more anthropomorphic device for better HCI; they might even help predict human actions in large computational models of human actors. We will explore some of these practical applications in § 4.

## 1.1 An Early Rebuttal of some Common Criticisms

One criticism made of these anthropomorphic formalisms is the argument that they don’t truly represent the trait they claim to. To address this point early, a responsibility formalism such as the one proposed need not be an entirely human-like representation

Hi Tim  
— this is  
what  
your  
todo  
notes  
look  
like! just  
use  
timnote.

Finish  
reading  
this!

of responsibility for every definition. Rather, there is a utility in an agent giving the *appearance* of responsibility. The utility is what is sought from creating these anthropomorphic algorithms for the effect it has on a system's behaviour — not perfect emulation of the trait itself. (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 it is useful to have computational frameworks for responsibility-like traits; we will see that these traits are indeed useful, and so that the criticism is moot. Computational trust formalisms are well documented as a valuable asset in solving HCI problems and designing aspects of intelligent agents, such as decision functions. We will see that computational responsibility follows in these footsteps, and has applications in AI and HCI just like trust. There are added benefits to responsibility formalisms, however, such as applications to a wider range of interdisciplinary study, and a very direct application in solving problems in areas like decision theory.

## 1.2 The Scope of the Model

Useful context for considering what sort of agents might be “responsible” can be found in exploring the agents to which our formalism may apply. As will be explored in § 3.4.1, there exist types of agents which we might not consider responsible in an ordinary setting, or for whom irresponsible behaviour might affect an assessment of responsible-ness. A human agent who is mentally handicapped, compared to a human agent with ordinary brain function but is lazy, wouldn't be seen as irresponsible when failing to hand an assignment in on time. To account for this difference in how agents' actions are accounted for in the proposed formalism, we limit the scope of what the formalism might model to agents who are:

- Reflective
- Interpretive
- Reactive

### 1.2.1 Reflective Agents

Sloman's work in the “space of minds” [Slo84] shows that artificial “minds” need not be remotely human-like. In order to limit the space of agents the proposed formalism would apply to, then, one might limit the space of minds those agents might inhabit.

There are several ways to limit the space of mind of the agents a formalism concerns: for example, Castelfranchi and Falcone [CF] define a “cognitive” agent as the lower limit of an agent's requirements for human traits for trust. They 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*.

A simpler way to state this, for the purposes of implementation in an artificially intelligent agent, would be that an intelligent agent should parametrise its decision function by its obligations. In this way, obligation to a certain goal or outcome influences actions chosen by that agent; this considering of obligation is required for those actions to be “responsible”, because an agent which does not parametrise by its obligations would have no way of accounting for what it ought to do when choosing an action by definition.

### 1.2.2 Interpretive, Reactive Agents

One can imagine other useful limitations of scope, too. For example, an agent should be able to interpret their own behaviours as responsible or irresponsible, such that they can assert the degree to which they should weight their obligations in their decision process:

Involved in an “interpretive 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, human-like subjectivity of responsibility can be simulated. We might go one step further, and more tightly constrain the subjective nature of an agent:

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 and other agents’ actions*.

Is the punctuation on the quotes fixed by inter-word spacing?

One can see that, with such a limitation on the agents a formalism concerns, the formalism becomes useful regardless of its computational application. While the formalism might be algorithmic in nature, the concepts behind it can be applied to social sciences also as the agents it concerns becomes more anthropomorphic. The interdisciplinary nature of a formalism such as this is a great asset in many areas, and allows for a common jargon when, for example, HCI researchers work with ethnographers in understanding the responsibility of a user.

The purpose of describing these terms is twofold. Partly, it is to introduce the notion that a proposed formalism of responsibility would be sensible in the types of agents it would target; this is useful to bear in mind when considering a trait which is normally human-specific developed as an algorithm. These terms are also introduced to show that the interdisciplinary jargon a formalism creates is naturally and easily defined by the formalism's construction. This shows in a concrete way that the formalism, even in an early form, has clear utility.

### 1.3 Proposal Overview

This proposal will be split into five main sections:

1. This introduction (§ 1), which lays the foundation for the research to be done and gives context for the background survey to come.  
This context is useful in framing the rest of the proposal with a useful perspective and an understanding of the field and its domain. This is particularly important when remembering that the formalism of social traits is not a familiar field to all computer scientists, unlike a field like algorithmic complexity.
2. A brief problem statement (§ 2), which details in specific terms the research intended to be undertaken.  
This will explore the research questions chosen, including the questions themselves, and why those questions are relevant and important to answer in the context that that introductory section gives.
3. A background survey (§ 3), which explores related literature to computational responsibility, including:
  - Mathematics and Social Sciences
  - Sociotechnical Systems research
  - Philosophical research

This will also complete the context provided in the introductory section by exploring components of non-computing science fields which impact the understanding of responsibility that the proposed work is founded upon, and the relevant work

Is the list of sections up to date?

Is the list of background topics up to date?

used in similar research projects for formalising social constructs. This will be used to inspire the direction the work takes.

4. A proposed approach (§ 4) to undertake the research suggested, which will explore potential options for the formalism and show how the relevant literature informs specific possible formalisms.  
Details of possible formalisms will be explored, and a course will be set for the direction the work will progress in.
5. A brief work plan (§ 5), which proposes a timeline for the work outlined in earlier sections, and shows in concrete terms the scope of the work and the steps taken to answer the research questions outlined in the earlier section.

## 2 Statement of Problem

With some background exploration on what a formalism of responsibility might entail, and an overview of its scope and utility, we can see that some formalism of responsibility has genuine utility. However, assessing how it might apply to artificial agents in practice requires the development of the formalism itself. It also requires that the formalism be applied to the specific category of minds outlined in § 1.2: reflective, interpretive, reactive agents.

We can address the feasibility of a real formalism which applies to these agents — and develop said formalism in the process — by answering 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?

I propose that this work would provide a valuable addition to the development of anthropomorphic trait formalisms, and that the work is also useful and interesting in its own right. These questions were also chosen in particular for a few reasons, which are detailed here.

### 2.1 Research by Scientific Method

A desired property of the questions chosen was that they should not be answerable by general insight from relevant literature. Instead, they should be answerable only by experimental method. This is because research questions which are answerable by literature insight are:



A: Best solved by a literature review, rather than scientific experiment, as where literature analysis is a viable research option it can be much more cost and time effective.

B: Inappropriate for masters level research and therefore unfit for the project at hand.

Therefore, research questions answerable only by experimental method are required, as experiment-driven research is the most appropriate research for a masters project, and questions answerable by literature review should in general be answered by this method (though there are of course exceptions).

The questions selected fit this criteria. This is because a responsibility formalism's efficacy in directing agent decisions and assisting an agent in assessing the responsibility of other agents is unclear upon the construction of the formalism. As will be explored, when the proposed formalism is fully defined, experiments must be undertaken in order to tweak the formalism's structure so as to effectively answer the questions — experimental data is required to determine whether agent decisions are suitably swayed and algorithms can properly assess other agents' responsibility. Therefore, the questions provided meet the requirement for experimental exploration.

## 2.2 Appropriately Fit Agent Context

In § 1.2, an argument was presented for certain traits agents must exhibit in order for applying a responsibility formalism to their behaviour to make sense. Appropriate research questions to be undertaken in the development of such a formalism should therefore limit the types of agents they apply to, such that those agents meet those criteria.

It can be seen that the research questions proposed limit the scope of the agents a formalism would apply to in just this way. For example, to direct the decisions made by an intelligent agent, that agent must make decisions taking the formalism into account — its decision function is therefore parametrised by the agent's responsibility. This makes that agent reflective.

For an agent to assess the actions of other agents, one would want that agent to be able to come to different conclusions about the actions of other agents than that agent did — it would therefore have to make use of a subjective interpretation function to come to these analyses. While this interpretation function isn't strictly necessitated by the agents a formalism would apply to — there are other ways to introduce subjectivity, such as differing training data on learning algorithms — it suggests development in this direction, as an interpretation function would be a very simple and not very limiting way to introduce this subjectivity. Parametrising this interpretation function by the current environment would also ensure an agent was reflective, making the interpretation function approach particularly appealing and admitting the final trait suggested by the problem context in § 1.2.

Given that the research questions proposed imply the correct types of agents as much as is possible, we will explore these questions in the literature that might help to indicate how to construct a solution to the problem. As experimental data will show whether a formalism will answer these research questions through an analysis of the “responsibleness” of their actions, relevant literature should inform the construction of the hypothetical formalism. This will in turn explore what an experiment which answers the proposed questions will look like.

### 3 Background Survey

Computational trust is a topic with a wealth of literature to draw inspiration from in an attempt to produce similar formalisms of responsibility. Marsh [Mar94] draws his own inspiration from literature as early as David Birkhoff’s 1930s work on creating an “Aesthetic Measure”, which was in effect a quantification of aesthetics. Much work on developing similar formalisms has been undertaken by a range of fields since.

The success of recent formalisms, such as Marsh’s, rest as all scientific work does on the research it serves to advance. As successful work has resulted from the study of social traits in social sciences and mathematics, a sensible approach to the construction of a new formalism would be to start from the same foundational concepts.

Unfortunately, not all of this related work will ultimately be of use in developing a formalism of responsibility — responsibility as a concept differs in important ways from trust. For example, trust is a concept often discussed in terms of an agent’s trust toward another: “Anne didn’t trust the contractor” describes Anne’s trusting relationship *toward* another actor. “The contractor performed their work responsibly”, however, is a description of the contractor’s responsibility as a trait — not as a relationship with Anne. Of course, this is one of a number of examples of the different natures of trust and responsibility. Fields which are not explored by trust might therefore be useful in elucidating some aspects of responsibility which trust-relevant literature may not cover. We find that these aspects are well covered by the field of Moral Responsibility, and similar philosophical work.

## 3.1 Social Sciences and Mathematics

### 3.1.1 Birkhoff's Aesthetic Measure

One of the earlier formalisms of a human factor<sup>12</sup> 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 (commonly "C&F"), as it is:

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

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.

---

<sup>1</sup>For 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.

<sup>2</sup>Also 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:

- 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.

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"[Deu62]:

- 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* if they choose to 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. Some examples are "Trust as Despair", "Trust as Virtue", "Trust as Masochism", and "Risk-taking or Gambling". 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, Luhmann 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... presupposes a situation of risk.”[Luh00]

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,  $x, y \in \mathcal{A}$ , 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 vacuum, 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,  $\alpha$ . 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 detailed 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 Birkhoff'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, FIRE, and others.

### 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 components to formulate a cognitive trust model:

1.  $x$ , a truster
2.  $y$ , a subject of trust
3.  $g$ , a goal of  $x$
4.  $\alpha$ , 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,  $\alpha$ . 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 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 encouraging 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 Eigentrust

Unlike the formalisms from Marsh and C&F, Eigentrust[KSGM03] is a trust formalism built without interdisciplinary work in mind. Instead, Eigentrust's main focus is that of software engineering: it is a formalism with the foremost principle of creating secure trust and reputation systems in a computer network. Indeed, the primary driving force of Eigentrust's mathematics is that of reputation, and it leverages linear algebra concepts together with a simple reputation system to create an effective trust formalism. Core to the reputation framework is that Eigentrust adds 1 to a score for a positive interaction, and -1 for a negative interaction, feeding these scores into trust scores which, through a series of matrices and eigenvalues, forms a global ledger of trust.

An intriguing feature of Eigentrust is that very global ledger of trust, shared in a peer-to-peer network. Having a shared knowledge of agents who are trustworthy or untrustworthy in a system acts as an interesting utility in exploring the value of a trust formalism. An example alluded to in the paper is that of downloads: if a download successfully completes from one peer to another, the receiving peer rates their server positively. Similarly, if the download fails, the receiving peer rates the server negatively. An agent seeking a certain download therefore attempts to interact with more trusted peers in the global ledger — this is a satisfying example of the value of a trust formalism in practical engineering.

In this way, Eigentrust creates a formalism that:

A: Represents gradations of trust

B: Allows for an effective distributed trust model

At first glance, Eigentrust appears to present a novel, elegant solution to trust formalism which might suit the responsibility formalism perfectly. However, nuances regarding its definition differ from that proposed, and significant work would be required to overhaul this model.

One reason for this is that Eigentrust explicitly does not attempt to interpret trust values; they are taken at face-value. As a result, the model wouldn't fit the space of minds of possible actors which was originally suggested as appropriate. The formalism constrains behaviour in other ways for appropriate trust modelling; however, some work may be required to produce a similar formalism for computational responsibility.

Another effect of adopting a framework like Eigentrust for building a responsibility formalism would be that it would work significantly less well for modelling an agent's personal obligations, making it less suitable for parametrising the decision function of an intelligent agent. An argument to be made is that Eigentrust presents a way for an agent to assess its environment. This is true. However, consider that a public ledger where one agent was seen to be more responsible than another would lead all agents to assign tasks to that supremely responsible agent; tasks would not therefore be assigned to marginally less responsible agents without some significant work to make the global scores from an Eigentrust-based model suitable for a responsibility model. Using Eigentrust-based scoring systems would also limit one's ability to separate types of responsibility in a similar way to Marsh's separation of basic, general and situational trust: Eigentrust's very simple  $+/-1$  reputation system permits little insight into types of actions an agent is adept at.

Eigentrust presents an interesting and unique approach to trust formalism which is practical and effective. Were the requirements of the responsibility formalism proposed any different, Eigentrust would posit a very enticing foundation on which the responsibility model could be based. Unfortunately however, its limitations pertaining to the more social form of responsibility desired, combined with its higher degree of mathematical complexity over alternatives, means that it is more suitable as a curiosity than a foundation for the purposes of this work.

### **3.3 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 responsibil-



ity modelling differs from a responsibility formalism in that a responsibility formalism needs to *generalise* reasoning 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 Sommerville’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 separation 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 separation of past responsibilities 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 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 § 3.4, 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.

## 3.4 Philosophy

Philosophy regarding moral responsibility is an area whose literature is both wide and deep. That said, not all moral responsibility literature is relevant to a computational responsibility 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.

### 3.4.1 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 concepts 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 inherent 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 social *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 be identified with, or indifferent to, it.

[Str62]

Strawson captures an essential component of Marsh's computational trust model: an actor's actions influence the perceived 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 Strawson’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<sup>+</sup>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 the actions for which an agent might consider itself responsible. 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 abnormal or 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 separation 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§ 1.2, 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.

### 3.4.2 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 Sommerville’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.

### 3.4.3 Deontic Logic

Deontic logic is a mathematical and philosophical logic of obligation.[vW51] It therefore would appear to lend itself perfectly to the task at hand; it would appear upon first glance that a responsibility formalism is already constructed. Unfortunately, this is not the case.

For one, deontic logic does deal with obligation, but from the perspective of formalising philosophical notions of action. Indeed, one of the insights a reader can expect to take from the work is a rigid definition of action. Deontic logic therefore deals with obliging to act as it contrasts to forbidding to act. Moreover, deontic logic creates relationships between different actions. For example:

$$drive \Rightarrow O(park\textit{safely})$$

is a deontic logic statement asserting that if one drives, one is obliged to park safely. However, the responsibility formalism required should assert a degree of responsibility for actions an agent might take at a given point in time — this is required to answer the first research question. As a result, deontic logic is a fundamentally inappropriate platform for the work at hand to be built upon.

Another issue is that deontic logic is a logical framework, rather than a mathematical formalism. As a result, it is difficult to model gradations of levels of responsibility from the perspective of deontic logic without making significant changes to it. Were the theory of computational responsibility proposed more similar to Castelfranchi & Falcone’s model than Marsh’s in its inspiration from computational trust formalisms, this would be less of an issue; however, Marsh’s approach to his formalism suits our philosophy regarding computational responsibility much better.

### 3.4.4 Sloman

While this is already explored somewhat in the introduction to give context to this proposal, Sloman's work on the space of possible minds[Slo84] is an interesting example of how intelligent agents might be addressed as increasingly anthropomorphic "minds". Whether an intelligent agent can have what is termed a mind is a philosophical question beyond the scope of this proposal. Sloman's work, however, does lend itself some useful concepts with which one might consider anthropomorphic agents.

Sloman's concept of the space of possible minds presents a notion that a mind, parametrised, can be represented by a series of spectra. These spectra combine as a mathematical Cartesian space, where a coordinate in that space represents a specific mind. There is, therefore, a subspace of this space which represents all possible human minds — a larger subspace which represents all possible biological minds — and possibly a broader still space which represents all possible minds, regardless of detail. Somewhere in this space, anthropomorphic algorithms reside.

The relevance and utility of Sloman's work to this responsibility formalism has already been presented: in acknowledging that an intelligent agent may not exist within the space of human minds, to impose responsibility upon the agent imposes human-like characteristics. In other words, the subspace of minds permissible by this formalism ought to be close to the subspace of human minds. To achieve this limitation, we impose upon the agents modelled by the formalism proposed traits of reflection, reaction, and interpretation.

## 3.5 Discussion

Computational Responsibility is a topic with a significant body of related literature, yet no specific literature on the topic. However, useful ideas expressed in the related literature — from both backgrounds of the humanities and the sciences — help to shine light on a path toward a formalism which will help to answer the research questions posited in § 2.

It's worth noting that Philosophy and Computing Science are unusual bedfellows. However, the deeply philosophical nature of this work shows that this need not be the case. Philosophical research may well have an impact on the ethics and cultural implications that advances in Computing Science make; it is important therefore to promote where possible interdisciplinary work of this nature. The philosophical literature reviewed is of great help in directing the project's next steps.

A formalism which does help to answer the research questions introduced will have to:

- Provide gradations of responsibility measures, in a similar way to that which Marsh's formalism permits

- Allow for the updating of the interpretation of responsibilities in a similar way to that which Strawson suggests
- Primarily address Sommerville’s “Causal” responsibilities, or (the given interpretation of) Scanlon’s “Attributive” responsibilities
- Follow Strawson’s note on how agents with human factors change their outlook with respect to these factors through interaction
- Account for responsibility on the level of personal responsibility and societal responsibility  
(This would be achieved by agents who act on individual responsibility but who as a collective use responsibility to lower negative exposure, as per Luhmann’s notes on risk)

In the proposed approach in the following section, we will explore possible formalisms of responsibility which leverage ideas from all of these useful insights from related literature, in order to better answer the research questions laid out. This will involve adopting a Marsh-like foundation for the formalism, and using interpretation functions and similar devices, with mathematical formalisms in the spirit of Birkhoff’s work on aesthetics. Its behaviour and approach to responsibility will be informed and directed by the relevant aspects of Strawson and Scanlon’s writing.

Eigentrust, as a foundation for this model, was unfortunately unsuitable. However, an open possibility is to adopt a similar reputation system for assessing the responsibility of other agents — this is unexplored in § 4, as it would not fit with the current conception of Marsh-like responsibility. Should such an approach be required, this will be assessed as the formalism is completed. The roadmap to completing the formalism is explored in § 5

### 3.5.1 The Relationship Between 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

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 re-appropriated (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 is the argument that task delegation is inherently trust; then, all one needs to 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 do also exist[LD08].

Develop arguments that the responsibility formalism might actually be put to good use, as per § 1

## 4 Proposed Approach

We can see that the literature available, while not on computational responsibility formalisms, is 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:

Identify risks in this approach

- The formalism should answer the research questions laid out in the problem statement§ 2.
- The formalism should suitably limit the scope of the actors it models to be reactive, reflective agents§ 1.2.  
The formalism doesn't apply to all agents; we can imagine some agents which are \*not\* responsible, such as those shown by Strawson§ 3.4.1, 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.
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 types of trust are useful for establishing the basic principles behind the responsibility model:

| <i>Marsh’s Trust</i> | <i>Responsibility Formalism</i> |
|----------------------|---------------------------------|
| Basic Trust          | Basic Responsibility            |
| General Trust        | General Responsibility          |
| Specific Trust       | Specific Responsibility         |

Table 1: The analogy between trust and responsibility through Marsh’s definitions

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:

| <i>Marsh's trust variables</i>   | <i>Variable Range</i> |
|--|-----------------------|
| Knowledge (of $x$ at time $t$ )  | True/False            |
| Importance (of knowing fact $x$ at time $t$ )                                    | $[0, +1]$             |
| Utility (of action $\alpha$ at time $t$ to an agent)                             | $[-1, +1]$            |
| Basic Trust (of agent $A$ at time $t$ )  | $[-1, +1]$            |
| General Trust (of agent $A$ at time $t$ in agent $B$ )                           | $[-1, +1]$            |
| Situational Trust (of agent $A$ at time $t$ in agent $B$ doing action $\alpha$ ) | $[-1, +1]$            |

Table 2: Overview of Marsh's model components

| <i>Proposed similar responsibility variables</i>  | <i>Proposed Variable Domains</i> |
|---|----------------------------------|
| Basic Responsibility (of agent $A$ at time $t$ )  | $[-1, +1]$                       |
| General Responsibility (of agent $A$ at time $t$ in agent $B$ )                           | $[-1, +1]$                       |
| Situational Responsibility (of agent $A$ at time $t$ in agent $B$ doing action $\alpha$ ) | $[-1, +1]$                       |

Table 3: Similar components a responsibility model can adopt in analogy to trust

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.

#### 4.2.2 Outlining the Proposed Formalism

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

| <i>Responsibility Term</i> | <i>Definition of Term</i>  |
|----------------------------|--|
| Agent / Actor              | Combination of Constraints, Environment, and Beliefs. 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 degree of obligation             |
| Action                     | Combination of Activity, Resource requirements, Possible Issues and Effect         |
| Authority                  | Agent which assigns an obligation to another agent                                 |

Table 4: Definitions of possible terms

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 of an obligation — these features change from goal to goal, and theoretically from agent to agent also. 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

My strategy for completing the formalism outlined is composed of four milestones, with concrete and well-defined deliverables:

1. Complete the formalism  
This will see a full formalism developed, complete with definitions.
2. Design a test case  
The product of this stage will be a scenario where the formalism can be fairly tested, so as to answer the research questions as laid out in § 4.3.
3. Test the formalism  
This will see intelligent agents developed and tested in the example scenario. It will produce experimental data, and may result in minor modifications to the formalism so as to properly simulate responsibility.
4. Write-up  
With experimental data collected and evaluated, the final report will be written up.

## 5.1 Completing the Formalism

While some detail regarding the formalism and its structure has already been surmised, the full extent of the formalism is incomplete.

One aspect of the formalism which has not been developed is the mathematical reasoning which turns the semantic descriptions of the formalism into one which is fully algorithmic. This will require some deep mathematical reasoning. It will also need to draw heavily on philosophical literature, so as to create a formalism which is socially and philosophically consistent. Other work yet to be completed involves the final definition of some components, which are a necessary part of the formalism but are unlikely to have a structural impact in the way that philosophical work might.

The formalism should be as close to the final product of the research as possible, so extra time is allotted to complete the formalism properly. This involves finalisation of the details, and checking against literature of other disciplines. It is intended that this be completed by mid-January.

## 5.2 Designing the Use Case

The use case required to test the formalism is particularly important, as the somewhat semantic nature of the research may make data collection and analysis difficult. It is believed though that a test case will arise naturally from both the related trust formalisms and their experimental technique, combined with the final formalism lending itself nicely to certain problems.

The test case should exhibit properties of the models specifically which answer the research questions proposed (§ 4.3). An AI strategy should also be devised during this stage. A nuance of this segment of work is that the experiment being developed must ultimately produce some data to analyse. Therefore, metrics which are appropriate for the analysis of an artificial agent's behaviour — and how responsibly that agent is behaving — should be produced.

As this segment of the work draws heavily on definitions from the formalism and previous formalisms of human traits, it is not expected to be a significant amount of work relative to others. It is expected to be completed around the beginning of February.

## 5.3 Testing the Formalism

Once the test case has been identified and necessary aspects of the experiment are established, the example scenario will be developed using agent modelling tools appropriate for the task at hand. The agents developed must be designed to make use of the responsibility formalism. This can be shown, for example, by observing that agents disregard

completely actions they are not responsible for. They also should identify and ignore other irresponsible agents. They should become more adept at this over time, in accordance with the research questions.

A concern in identifying time required for this block of work is that the formalism may prove infeasible computationally once implementation is attempted, requiring further refinement. To account for this, the analysis of the data produced should end this block of work around the middle of March.

## 5.4 Write-up

Once experimental data has been collected, verified, and analysed, the final block of work is to finish the write-up of the project. This segment will include writeup of analysis and production of visualisations, as well as incorporating any writing produced as the project progresses. It is expected that, should the rest of the work be properly completed at this point, that the report be finished around the beginning of April.

## 5.5 Review of Time Allocation

The time allocated to the various components of the proposed work is therefore:

| <i>Section</i>             | <i>Intended completion time</i> |
|----------------------------|---------------------------------|
| Finalise Formalism         | Mid-January                     |
| Design Use Case            | Beginning of February           |
| Run Tests and Collect Data | Early / Mid-March               |
| Write-up Completed         | End March                       |

Table 5: Time allocated to various parts of the work

Note that the completion of the work ends three weeks before the project is due. To avoid complications and rushed research, the work has been planned with a three-week margin of time should any part overrun. This is due to the change that revisions of the formalism may be required, and iterating the work may be time consuming and difficult to account for. In addition, providing a small margin of time ensures minimal rushing of the project and appropriate time management indicates more quality overall work.

## References

- [Bir] G D Birkhoff. AESTHETIC MEASURE.
- [CE11] Partheeban Chandrasekaran and Babak Esfandiari. A model for a testbed for evaluating reputation systems. *Proc. 10th IEEE Int. Conf. on Trust, Security and Privacy in Computing and Communications, TrustCom 2011, 8th IEEE Int. Conf. on Embedded Software and Systems, ICESS 2011, 6th Int. Conf. on FCST 2011*, pages 296–303, 2011.
- [CF] Cristiano Castelfranchi and Rino Falcone. Social Trust: A Cognitive Approach.
- [Deu62] Morton Deutsch. Cooperation and trust: Some theoretical notes. 1962.
- [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.
- [HJS06] Trung Dong Huynh, Nicholas R. Jennings, and Nigel R. Shadbolt. An integrated trust and reputation model for open multi-agent systems. *Autonomous Agents and Multi-Agent Systems*, 2006.
- [Hon] Ted Honderich. Free will, determinism and moral responsibility – the whole thing in brief.
- [KSGM03] Sepandar D Kamvar, Mario T Schlosser, and Hector Garcia-Molina. The eigentrust algorithm for reputation management in p2p networks. In *Proceedings of the 12th international conference on World Wide Web*, pages 640–651. ACM, 2003.
- [LD08] Emiliano Lorini and Robert Demolombe. From binary trust to graded trust in information sources: a logical perspective. In *International Workshop on Trust in Agent Societies*, pages 205–225. Springer, 2008.
- [LSSB10] Russell Lock, Tim Storer, Ian Sommerville, and Gordon Baxter. Responsibility modelling for risk analysis. 2010.
- [Luh00] Niklas Luhmann. Familiarity, confidence, trust: Problems and alternatives. 2000.
- [Mar94] Stephen Paul Marsh. Formalising Trust as a Computational Concept. *Computing*, Doctor of(April):184, 1994.
- [MBEK<sup>+</sup>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.

- [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.
- [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.
- [vW51] G. H. von Wright. Deontic logic. *Mind*, 60(237):1–15, 1951.