# Agent-Based Logics in Dependent Type Theory

Colm Baston

Thesis submitted to the University of Nottingham for the degree of Doctor of Philosophy, September 2024



# Abstract

This thesis is on the formalisation of mathematics in Martin-Löf type theory. This is a class of dependently-typed functional programming languages whose rules form a language suitable for the statement and proof of mathematical theorems. Programs in type theory can be mechanically checked to be well-typed, which corresponds to proofs being checked to be valid. In particular, we explore the semantics of certain agent-based logics, and seek to give them appropriate representations using the concepts of type theory.

We embed the syntax of epistemic modal logic in type theory as predicates over a type of possible worlds. Knowledge operators are defined in this setting by stating a set of properties they must satisfy. We prove this knowledge operator semantics equivalent to the traditional relational semantics of epistemic logic. Common knowledge is then defined in the embedding as a coinductive predicate, whose proofs are infinite data structures. We prove this representation is equivalent to the intuitive definition as iterated universal knowledge, and prove it is equivalent to the relational interpretation. In coalition logic, we represent game forms and playable effectivity functions in type theory, and outline a proof of their equivalence.

# Acknowledgements

Foremost, I would like to thank Venanzio Capretta, who agreed to supervise this project and whose knowledge and research expertise have been indispensible ever since. I am also very appreciative of Graham Hutton, who has been very encouraging from the start, but took on the role of my second supervisor midway through, and was of great assistance towards the end of my studies.

More broadly, I am grateful to the past and present members of the Functional Programming Lab, and others from the School of Computer Science, for many insightful discussions over the years. In addition, I greatly enjoyed the undergraduate teaching opportunities offered by the school, which gave a welcome break from research.

Thanks to Nicolai Kraus and Roy Crole for their constructive comments and suggestions in finalising this thesis. Finally, without the support of my parents, Peter and Noëlle, completion would not have been possible.

# Contents

Abstract     i       Acknowledgements     i							
					Ι	In	Introduction
	1.1	Thesis Structure and Contributions	2				
II	De	ependent Type Theory	5				
	2.1	Background	5				
	2.2	Simple Types	6				
	2.3	Type Theory as a Logic	13				
	2.4	Dependent Types	16				
	2.5	Equality and Equivalence	21				
	2.6	Inductive Types	28				
	2.7	Coinductive Types	33				
II	ΙE	Distemic Modal Logic	41				
	3.1	Background	41				
	3.2	Epistemic Logic and Relational Semantics	43				
	3.3	Type-Theoretic Embedding	45				
	3.4	Knowledge Operator Semantics	49				
	3.5	Correspondence with Relational Semantics $\ldots \ldots \ldots \ldots$	51				
IV	Co	oinductive Common Knowledge	57				
	4.1	Background	57				
	4.2	Coinductive Definition	62				
	4.3	Relational Definition	65				
	4.4	Coinductive and Relational Equivalence	67				
$\mathbf{V}$	Co	palition Logic	71				
	5.1	Background	71				
	5.2	Game Form Construction	73				

5.3	Decidable Subsets	75			
5.4	Partial Functions	78			
5.5	Future Work	80			
VI C	oncluding Remarks	83			
A A	gda Formalisations	85			
A.1	Dependent Type Theory	85			
A.2	2 Epistemic Modal Logic	95			
Α.	3 Coinductive Common Knowledge	98			
Bibliography 10					

# Chapter I Introduction

Martin-Löf type theories are a class of functional programming languages which give strict rules for the construction and use of different *types* of data. For example, a function defined to take a Boolean value as argument cannot be applied to a natural number, and vice versa. While languages with strong type systems are common throughout computer science and software engineering, these most often treat data types and their values as completely separate entities that cannot be mixed. In type theories, however, types are themselves first-class values, able to be stored in data structures, returned by functions, and so on, giving rise to *dependent types* — types which depend on values.

Dependent types have applications in common programming tasks. A type for matrices may depend on natural number values which state their dimensions. A matrix multiplication function can then be defined which only multiplies  $a \times b$ and  $b \times c$  matrices, with a matching dimension b, and the type of the output matrix is also labelled with its size,  $a \times c$ . The type system ensures statically, at the time the program is compiled, that all matrix multiplications are applied to appropriately-sized matrices, with no runtime checks necessary.

Simultaneously, through the phenomenon of the Curry-Howard correspondence, Martin-Löf type theory serves as a foundation for mathematics. Types can be seen as mathematical statements, their values as proofs for the statement, with the specific type giving the valid ways of constructing a valid proof of the statement. Dependent types are predicates whose truth depends on the values to which they are applied, and type theory has connectives analogous to the universal and existential quantifiers of a predicate calculus.

There are many implementations of dependently-typed languages, some of which are specifically designed as *proof assistants*, with features that assist in the construction of mathematical proofs. The type-checking apparatus of their compilers double as proof checkers — if a program compiles without error, the corresponding mathematical proofs are validated to be correct. Even within computer science, mathematics has traditionally been done on pen and paper, new theorems communicated in the literature through prose. While care is taken to ensure constructions are compatible with a mathematical foundation, most often a system of set theory, rarely are proofs rigorously checked by a mechanical process, and mistakes sometimes slip through. In recent years, even pure mathematicians with no previous connection to computer science have turned to type theory as a potential solution to this problem [8, 11].

This thesis is on formalising mathematical results in the language of type theory, with the objective that they can be rigorously verified by a proof assistant. In particular, it explores the semantics of certain agent-based logics, epistemic modal logic and coalition logic, and their representation within type theory.

In epistemic modal logic, an embedding of the logic is defined in the language of type theory, with proof that the embedding corresponds to the traditional semantics. This allows features of the type theory to be used in epistemic logic, for example, the notion of common knowledge is be defined using the feature of coinductive types. In coalition logic, we focus on the formalisation of an equivalence proof between two semantic interpretations of the logic.

#### **1.1** Thesis Structure and Contributions

The thesis aims to be accessible to those with a mathematical background, but with no prior exposure to Martin-Löf type theory, so Chapter II begins with a brief exposition. This is established material which can be pieced together from other sources, but the chapter provides a consistent syntax and an introduction to concepts that will be essential in later chapters. This chapter has a corresponding formalisation in the Agda proof assistant, provided in Appendix A.1, whose definitions are then used as a basis for later formalisations.

The content of Chapters III and IV is based on the following publication [5].

 THE COINDUCTIVE FORMULATION OF COMMON KNOWLEDGE Joint work with Venanzio Capretta Published in the Proceedings of the 9<sup>th</sup> International Conference on Interactive Theorem Proving, Oxford 2018

In Chapter III, we consider *epistemic modal logic*, which adds operators to the syntax of a propositional logic which model the knowledge of agents. The chapter begins with an introduction to the syntax of the logic and its traditional semantics, based on equivalence relations, and goes on to develop a representation for them in type theory.

This contributions of this chapter are as follows.

#### 1.1. THESIS STRUCTURE AND CONTRIBUTIONS

- An embedding for the syntax of the logic is defined in type theory, based on predicates over a type of possible worlds, Definitions 3.3.2 and 3.3.3.
- A semantic specification for knowledge operators within this embedding, Definition 3.4.2. In particular, this knowledge operator semantics introduces a new, infinitary deduction rule, *preservation of semantic entailment*, Definition 3.4.1, which is proved to subsume two axioms of epistemic logic, Theorems 3.4.3 and 3.4.4.
- Transformation functions between the knowledge operator semantics and the relational semantics, Definitions 3.5.2 and 3.5.4, with proof that the transformations yield the opposite semantics, Theorems 3.5.3 and 3.5.9.
- Proof that these transformations are inverse, forming an isomorphism up to propositional equivalence, Theorems 3.5.10 and 3.5.11.
- A formalisation of these results in the Agda proof assistant, Appendix A.2.

Chapter IV continues the focus on epistemic modal logic, extending it with multiple agents, each with their own knowledge operator. The concept of *common knowledge* is introduced, which exists among groups of agents when a fact is known universally, the fact of universal knowledge is itself known universally, and so on ad infinitum. The semantics of the common knowledge operator is explored in the context of a famous logic puzzle, which involves common knowledge in a subtle way, before going on to define common knowledge in type theory.

The contributions of this chapter are as follows.

- In the setting of the embedding of Chapter III, an operator for common knowledge is defined as a coinductive type, Definition 4.2.5, and is shown to satisfy the knowledge operator properties, Theorem 4.4.7.
- The coinductive common knowledge operator is shown to correspond to the intuitive idea of iterated universal knowledge, Theorems 4.2.10 and 4.2.11.
- The relational interpretation of common knowledge is transformed, using the isomorphism of the previous chapter, into a common knowledge operator, Definition 4.3.2.
- The coinductive and relational versions of the common knowledge operator are proven to coincide, Theorem 4.4.2.
- A formalisation of these results in the Agda proof assistant, Appendix A.3.

The content of Chapter V is based on the following publication [6].

 GAME FORMS FOR COALITION EFFECTIVITY FUNCTIONS Joint work with Venanzio Capretta Published in the Proceedings of the 25<sup>th</sup> International Conference on Types for Proofs and Programs, Oslo 2019

In Chapter V, we turn our attention to *coalition logic*, in which agents can act together towards a common set of goals. Coalition logic has two equivalent semantics based on game forms and playable effectivity functions, and an outline of the equivalence proof is given, with the objective of a formalisation of the proof in type theory.

The proof is involved and heavily set-theoretic, posing significant challenges for formalisation in type theory. This was an ambitious aim late into the research period, and the formalisation is not complete enough to be fully-verified by proof assistant. However the chapter proposes solutions to some of these challenges. These include representation of subsets with decidable membership as Boolean functions, Definition 5.3.1, and the use of partial functions to define type families over only a subset of a type at a time, Definition 5.4.2.

# Chapter II Dependent Type Theory

This chapter introduces a formal system of dependent type theory which will be used for the developments in later chapters. The coverage aims to be sufficient to gain an intuition, but not completely comprehensive — it won't be possible to implement a type-theoretic proof assistant based on this chapter alone.

The basic rules of type formation, introduction, and elimination are discussed in the context of a simply-typed  $\lambda$ -calculus. We then go on to say how this system forms a mathematical foundation via the Curry-Howard correspondence and how it differs from systems based on set theory. Dependent types and propositional equality are introduced, giving rise to a system of intuitionistic higher-order logic. Eventually, we will cover inductive types, which allow for the definition of wellfounded tree-like structures in the type theory, and coinductive types, which are non-wellfounded trees with potentially infinite depth.

## 2.1 Background

In the late 19<sup>th</sup> and into the early 20<sup>th</sup> century, significant effort was being made to base the previous discoveries of mathematics in a unified foundation of mathematical logic. Foremost among the proposed systems was a theory of sets, where all mathematical structures could be built up from sets, with formal rules for set formation and to determine set membership. However, it was discovered that any such system which includes an unrestricted comprehension principle is inconsistent, a result today known as Russell's paradox — does the set of all sets that are not members of themselves contain itself?

Bertrand Russell proposed his theory of types to avoid the paradox [53]. There is a hierarchy of *types*, with each set being assigned to unique type. Instead of adopting an unrestricted comprehension principle, sets may only be composed of elements of a strictly lower type in the hierarchy. This stratification of types renders the problematic Russell set inexpressable. The primary system of concern to this thesis is the type theory due to Per Martin-Löf [43,44]. While it may bear little resemblance to Russell's theory, the idea of types remains an essential concept. Each object is assigned a unique type, and operations may only be performed on objects of appropriate types. While this may seem restrictive when compared to more free-form set theories, a welldesigned type system ultimately provides a great amount of utility, as well as desirable computational properties.

Based on Alonzo Church's typed  $\lambda$ -calculus [16] and its later extensions, such as System F [29, 50], Martin-Löf type theory serves as both a dependently-typed functional programming language and a foundation for mathematics. It may be used to implement algorithms, but also to state and prove mathematical statements, all within the same language. Throughout this thesis, "dependent type theory", or simply "type theory", will refer specifically to a system of Martin-Löf type theory.

## 2.2 Simple Types

The constructions of type theory, called *terms*, are assigned to a unique type from the moment they are introduced and only exist in the context of their type. The language of type theory is defined by *judgements*, the most fundamental of which is the typing relation.

$$\Gamma \vdash t : T$$
 "in context  $\Gamma$ , term t has type T"

The context  $\Gamma$  is included here so that typing judgements may follow from prior assumptions. For the purposes of this introduction, a context may be considered a mapping of variable names to their types:  $v_1 : T_1, \ldots, v_n : T_n$ . When a judgement is asserted unconditionally, with an empty context of assumptions, it is common to omit the context in the notation. For example, when we introduce the type of natural numbers, we may simply write **zero** :  $\mathbb{N}$ .

Before defining other types, we fix a *type universe* Type — a type whose terms are themselves types. The judgement T: Type is understood to mean that T is a type as well as a term, and it may in turn have its own terms. If we derive  $\Gamma \vdash t : T$  without having first established  $\Gamma \vdash T :$  Type, this is considered to be an implicit premise.

We define the valid judgements of a type theory using inference rules in the style of natural deduction. Our first types, the empty type 0, and the unit type 1 containing exactly one term, are defined axiomatically.

 $\bigcirc -\text{FORM} \quad \boxed{\square : \text{Type}} \qquad \boxed{\square -\text{FORM} \quad \boxed{\square : \text{Type}}} \qquad \boxed{\square -\text{INTRO} \quad \cancel{\star : \square}}$ 

Rules for creating new types are called *formation rules*, while rules for creating new terms of a type are called *introduction rules*. We give a formation rule for each of the new types, but only 1 has an introduction rule to define its term  $\star$ , as 0 is empty, uninhabited by terms.

At this point, we could go on to define the Boolean type 2, containing two terms in a similar way, with one formation rule and two introduction rules, but it is more convenient to define it in terms of a compound type. We define a binary sum operator on types that forms a new type combining both types' terms, tagged with which of the summed types the term originated.

$$+-\text{FORM} \frac{\Gamma \vdash A : \text{Type} \qquad \Gamma \vdash B : \text{Type}}{\Gamma \vdash A + B : \text{Type}}$$
$$+-\text{INTRO}_{\text{inl}} \frac{\Gamma \vdash a : A}{\Gamma \vdash \text{inl} \ a : A + B} \qquad +-\text{INTRO}_{\text{inr}} \frac{\Gamma \vdash b : B}{\Gamma \vdash \text{inr} \ b : A + B}$$

With this type formation rule, without considering nesting, we can form four new types with each combination of  $\mathbb{O}$  and  $\mathbb{1}$ . Since neither side of  $\mathbb{O} + \mathbb{O}$  has any way of constructing a term, this type is uninhabited, but the other types have the following terms.

$$\mathsf{inl} \star : \mathbb{1} + \mathbb{0}$$
  $\mathsf{inr} \star : \mathbb{0} + \mathbb{1}$   $\mathsf{inl} \star : \mathbb{1} + \mathbb{1}$   $\mathsf{inr} \star : \mathbb{1} + \mathbb{1}$ 

We will use the type 1 + 1 as our Boolean type, with the two terms listed above representing true and false. We introduce a new judgement to state when two terms are definitionally equal.

```
\Gamma \vdash a = b : T "in context \Gamma, a and b are equal terms of type T"
```

To be useful as a notion of equality, we need basic rules to endow some standard properties on this relation: reflexivity, symmetry, transitivity, and congruence rules, omitted here. These rules ensure that the type theory makes no distinction between the two terms — one can always be used in place of the other.

For convenience, we can add rules to assert 2 = 1 + 1, true = inl  $\star$ , and false = inr  $\star$ . It is important to note that there is no intrinsic value in these names, they are just syntactic identifiers — we could just as easily define true and false the other way around.

The choice of 1+1 to represent the Boolean type is appropriate because it has exactly two terms. That is, it can be shown that any term of the type constructed in an empty context will be judgementally equal to either inl  $\star$  or inr  $\star$ , and these terms are distinct from one another. This is unsurprising at the moment, as the only way of constructing terms for 1+1 is by using the introduction rules directly, but will remain to be the case as new rules are added to the system. We will later be able to observe this fact from within the type theory, proving it as a theorem using an internal notion of propositional equality introduced in Section 2.5.

A second compound type is the product type. It is defined as another binary operator on types which results in the type of pairs formed by terms of the corresponding types.

$$\times \operatorname{-Form} \frac{\Gamma \vdash A : \mathsf{Type} \qquad \Gamma \vdash B : \mathsf{Type}}{\Gamma \vdash A \times B : \mathsf{Type}}$$
$$\times \operatorname{-Intro} \frac{\Gamma \vdash a : A \qquad \Gamma \vdash b : B}{\Gamma \vdash \langle a, b \rangle : A \times B}$$

As an example, the type  $2 \times 1$  consists of pairs where the first component is a Boolean and the second is  $\star$  — it has only two terms  $\langle \text{true}, \star \rangle$  and  $\langle \text{false}, \star \rangle$  that can be distinguished by judgemental equality. Similarly,  $2 \times 2$  has four terms  $\langle \text{true}, \text{true} \rangle$ ,  $\langle \text{false}, \text{true} \rangle$ ,  $\langle \text{true}, \text{false} \rangle$ , and  $\langle \text{false}, \text{false} \rangle$ , while any product of  $\mathbb{O}$ is uninhabited since no terms may be constructed to occupy the corresponding position in the pair.

As the notation suggests, the sum and product types parallel sums and products of arithmetic, and we can form types with any desired, finite number of terms by nesting them. For example, 1 + (1 + 1) is a type of exactly three terms, up to judgemental equality. This analogy holds when considering the properties of associativity, commutativity, and distributivity — the type (1 + 1) + 1 also has three terms. When including sums of  $\mathbb{O}$ , there are infinitely many ways to form *equivalent* versions of each finite type.

These different constructions, though equivalent, are considered distinct types with distinct terms — the judgement 1 + (1 + 1) = (1 + 1) + 1 cannot be derived, the term inr (inl \*) has the left type but not the right, and so on. For this reason, the type theory presented here is called *intensional*, as opposed to *extensional*. This distinction will be revisited, and a precise definition of type equivalence given, alongside the introduction of propositional equality in Section 2.5.

We've seen how to construct some basic types, but we can do little else with them at present. The function type is at the core of how we introduce computation into the system. Like sums and products, it is a binary operator on types.

$$\rightarrow \text{-Form} \frac{\Gamma \vdash A : \text{Type} \quad \Gamma \vdash B : \text{Type}}{\Gamma \vdash A \rightarrow B : \text{Type}}$$
$$\rightarrow \text{-Intro} \frac{\Gamma, v : A \vdash b : B}{\Gamma \vdash \lambda(v : A). \ b : A \rightarrow B}$$

The premise of the introduction rule shows the first non-trivial usage of contexts —  $\Gamma, v : A$  denotes the context formed by extending  $\Gamma$  with an additional assumption. Term b is allowed to have occurrences of variable v within its substructure, so the judgement b: B is conditional on the assumption of v: A. In the conclusion of the rule, the assumption is discharged and term  $\lambda(v: A)$ . b is now a function of variable v, resulting in term b. For notational convenience, we may sometimes drop the type of the bound variable where it can be inferred and would provide no additional clarity, simply writing  $\lambda v. b$ .

Strictly speaking, all functions in type theory are unary, taking exactly one input and producing exactly one output. Functions of higher arity may be simulated using products or, as is more common in functional programming, by *Currying* [21] — returning another function which is responsible for handling the next argument. For example, the types  $A \times B \to C$  and  $A \to (B \to C)$  both represent binary functions, the latter being a Curried version of the former.

To apply a function to an argument, we add an *elimination rule* to the type theory. Elimination rules specify how terms of a particular type may be used in computation, complementing the introduction rules for that type.

$$\rightarrow \text{-ELIM} \frac{\Gamma \vdash f : A \rightarrow B \quad \Gamma \vdash a : A}{\Gamma \vdash f \; a : B}$$

Application of the function f to argument a is denoted here by juxtaposition in the tradition of the  $\lambda$ -calculus and functional programming. The actual result of the function application is specified as a judgemental equality, the  $\beta$ -reduction rule.

$$\beta \text{-Reduce} \ \frac{\Gamma \vdash (\lambda(v:A).\ b): A \to B \qquad \Gamma \vdash a: A}{\Gamma \vdash (\lambda(v:A).\ b)\ a = b[v/a]: B}$$

In the rule's conclusion, b[v/a] denotes syntactic substitution of all free occurrences of the variable v by term a inside the substructure of b. We shall omit the full intricacies of  $\alpha$ -conversion and variable capture here. It is sufficient to say that the choice of variable name is arbitrary, and we identify  $\lambda$ -terms which only differ by the names of their bound variables. If free variables in a would be captured by naïve substitution of b[v/a], we first perform  $\alpha$ -conversion — all free occurrences of v in b are renamed such that they are not free in a.

An  $\eta$ -reduction rule is also defined. This states that a  $\lambda$ -term which only passes its argument through to another function may be identified with that function.

$$\eta\text{-REDUCE} \frac{\Gamma \vdash f : A \to B}{\Gamma \vdash (\lambda(v : A). f v) = f : A \to B}$$

Furthering the arithmetic analogy, function types are exponentials. Interpreting types A and B as cardinal numbers, the function type  $A \to B$  is inhabited by  $B^A$  functions that can be distinguished pointwise. For example, there are no functions  $1 \to 0$ , one function  $0 \to 1$ , four functions  $2 \to 2$  — the identity, a constant function for each of **true** and **false**, and Boolean negation — and so on. The rules given so far allow us to define some of these Boolean functions.

$$\begin{aligned} \mathsf{id}_2 &= \lambda b. \ b &: 2 \to 2\\ \mathsf{const}_\mathsf{true} &= \lambda b. \ \mathsf{true} \ : 2 \to 2\\ \mathsf{const}_\mathsf{false} &= \lambda b. \ \mathsf{false} : 2 \to 2 \end{aligned}$$

We cannot yet define Boolean negation as we can only interact with the function argument generically — we cannot do one thing when the argument is true and a different thing when it is false. To remedy this, we add elimination and reduction rules to our basic types. First, we give elimination rules for 0 and 1.

$$\mathbb{O}\text{-}\text{ELIM} \frac{\Gamma \vdash e : \mathbb{O}}{\Gamma \vdash \mathsf{elim}_{\mathbb{O}} \ e : T} \qquad \mathbb{1}\text{-}\text{ELIM} \frac{\Gamma \vdash t : T \qquad \Gamma \vdash u : \mathbb{1}}{\Gamma \vdash \mathsf{elim}_{\mathbb{O}} \ t \ u : T}$$

In each case, these say how to transform an element of the type into a term of an arbitrary type T. For 1, we must first be given a t:T, but for  $\mathbb{O}$  we allow this without restriction as it should not be possible to construct an  $e:\mathbb{O}$  in an empty context. If we are able to derive  $e:\mathbb{O}$ , it is only because it is implied by  $\Gamma$ , such as in the case of a function whose domain is  $\mathbb{O}$ . To illustrate this, we give a derivation for the function  $\lambda e. \operatorname{elim}_{\mathbb{O}} e:\mathbb{O} \to \mathbb{1}$  which uses this rule.

$$\begin{array}{c}
\mathbb{1} \text{-Form} \overline{1: \text{Type}} \\
\mathbb{0} \text{-ELIM} \overline{e: 0 \vdash 1: \text{Type}} \quad \overline{e: 0 \vdash e: 0} \\
\rightarrow \text{-INTRO} \overline{\frac{e: 0 \vdash \text{elim}_0 \ e: 1}{\lambda e. \ \text{elim}_0 \ e: 0 \rightarrow 1}}
\end{array}$$

This derivation uses two unlabelled rules which haven't been discussed. One is context weakening, the other an identity axiom for when a concluded typing judgement is already assumed in the context.

For the reasons outlined above, there is no reduction rule for  $\mathbb{O}$ . We need to give one rule for for 1 as it has only one term, simply reducing to the provided value t: T when applied  $\star$ . The rule therefore does no meaningful work, simply acting as a constant function returning t.

$$1-\text{Reduce} \frac{\Gamma \vdash t:T}{\Gamma \vdash \mathsf{elim}_1 \ t \ \star = t:T}$$

The sum type A + B has two cases, one when the term is of the form inl a, carrying term a : A, and the other for inr b, where b : B. The elimination and reduction rules describe how we can compute a value of type T in either case.

$$\begin{aligned} +-\text{ELIM} & \frac{\Gamma \vdash f: A \to T \quad \Gamma \vdash g: B \to T \quad \Gamma \vdash s: A + B}{\Gamma \vdash \mathsf{elim}_{+} f \ g \ s: T} \\ +-\text{REDUCE}_{\mathsf{inl}} & \frac{\Gamma \vdash f: A \to T \quad \Gamma \vdash g: B \to T \quad \Gamma \vdash a: A}{\Gamma \vdash \mathsf{elim}_{+} f \ g \ (\mathsf{inl} \ a) = f \ a: T} \\ +-\text{REDUCE}_{\mathsf{inr}} & \frac{\Gamma \vdash f: A \to T \quad \Gamma \vdash g: B \to T \quad \Gamma \vdash b: B}{\Gamma \vdash \mathsf{elim}_{+} f \ g \ (\mathsf{inr} \ b) = g \ b: T} \end{aligned}$$

To compute a value of type T from A + B, we provide a function f to handle the inl case and a function g to handle inr. One reduction rule is specified for each case, applying the appropriate function to the carried term.

Recall that we defined 2 = 1 + 1, true = inl  $\star$ , and false = inr  $\star$ . We can therefore use  $\mathsf{elim}_+$  and  $\mathsf{elim}_1$  to define functions which can discriminate the Boolean value passed to them, such as Boolean negation.

$$\mathsf{not}_2 = \mathsf{elim}_+ \ (\mathsf{elim}_1 \ (\mathsf{inr} \ \star)) \ (\mathsf{elim}_1 \ (\mathsf{inl} \ \star)) : 2 \to 2$$

The first function will handle the true case, converting inl into inr, and vice versa for the second function which handles false. Note that using the  $\eta$ -reduction rule, elim<sub>+</sub> is treated here as a Curried function, being partially-applied to only two of its arguments, leaving the final Boolean argument implicit.

The terms of the product type  $A \times B$  may only be of the form  $\langle a, b \rangle$ , where a : A and b : B. We therefore provide the eliminator  $\operatorname{elim}_{\times}$  with a single Curried function to handle both a and b at once.

$$\begin{split} \times -\text{ELIM} & \frac{\Gamma \vdash f: A \to B \to T \quad \Gamma \vdash p: A \times B}{\Gamma \vdash \mathsf{elim}_{\times} \ f \ p: T} \\ \times -\text{Reduce} & \frac{\Gamma \vdash f: A \to B \to T \quad \Gamma \vdash a: A \quad \Gamma \vdash b: B}{\Gamma \vdash \mathsf{elim}_{\times} \ f \ \langle a, b \rangle = f \ a \ b: T} \end{split}$$

For example, we can now define projection functions to extract just one part of the pair, or a function to swap the components of a pair. These examples assume a context  $\Gamma$  including A: Type and B: Type.

$$\begin{array}{ll} \Gamma \vdash & \pi_1 = \mathsf{elim}_{\times} \ (\lambda a. \ \lambda b. \ a) & : A \times B \to A \\ \Gamma \vdash & \pi_2 = \mathsf{elim}_{\times} \ (\lambda a. \ \lambda b. \ b) & : A \times B \to B \\ \Gamma \vdash \mathsf{swap} = \mathsf{elim}_{\times} \ (\lambda a. \ \lambda b. \ \langle b, \ a \rangle) : A \times B \to B \times A \end{array}$$

The rules introduced up to now constitute a simply-typed  $\lambda$ -calculus extended with sums and products. We will now observe some important computational properties exhibited by the system. The judgemental equality relation presented here is symmetric — it can be applied in either direction. However, it is useful to consider a one-way reduction relation which simplifies the term at each step in order to compute a final result. We can form such a relation by taking the reduction rules of the system and substituting = by  $\rightsquigarrow$ . In each rule, the equation has been arranged so the leftto-right direction represents a simplification step.

We don't imbue  $\rightsquigarrow$  with the equivalence relation properties of equality but do allow congruence rules, so subterms may be reduced without affecting outer parts of an expression. We also define  $\rightsquigarrow^*$  as the reflexive, transitive closure, representing a sequence containing any number of individual reduction steps. A term with no possible reductions is called a *normal form*, and the process of applying reductions to reach a normal form is called *normalisation*.

For example, take the example expression swap  $\langle not_2 \text{ true}, \star \rangle : \mathbb{1} \times \mathbb{2}$  which uses some of our previous definitions. Here are two reductions for this term, one first reducing swap, the other first reducing  $not_2$ .

	swap $\langle not_2  true, \star  angle$		swap $\langle not_2  true, \star  angle$
$\rightsquigarrow^*$	$\langle \star,  not_2  true  angle$	$\rightsquigarrow^*$	swap $\langle false, \star \rangle$
$\rightsquigarrow^*$	$\langle \star,  false \rangle$	$\rightsquigarrow^*$	$\langle \star,  false \rangle$
=	$\langle\star, inr \star \rangle$	=	$\langle\star,inr\star angle$

This illustrates the property of *confluence* — when two different reduction steps are possible from the same expression, there will always be subsequent reductions which lead both forks back together. The choice of which reduction to use gives rise to reduction strategies such as normal-order and applicative-order reduction, and the concept of lazy evaluation as is implemented in the functional programming language Haskell [42].

Furthermore, the system is strongly normalising — the reduction relation is wellfounded, so a normal form will always be reached in a finite number of steps. This is a consequence of only allowing well-typed terms, which prohibits divergent terms that are possible in the untyped  $\lambda$ -calculus, such as the term  $(\lambda\omega. \omega \omega) (\lambda\omega. \omega \omega)$ , which reduces to itself in a single step. In a typed setting, however, the self-application  $\omega \omega$  is forbidden as it cannot be assigned a type.

Additionally, normal forms exhibit *canonicity* — for each type former excluding functions, closed normal forms are of canonical form. That is, when a term can be constructed in an empty context, its normal form will have one of the syntactic patterns defined by the introduction rules for its type, called its constructors. Since  $\mathbb{O}$  has no introduction rules, and therefore no constructors, it also has no canonical forms, while 1 has a single constructor and canonical form  $\star$ . A sum's normal forms will always be constructed by one of inl or inr, and a product's by the pairing constructor  $\langle \_, \_ \rangle$  — the subterms carried by the constructors will themselves be of canonical form with respect to their own types.

Finally, the typing and equality judgements are *decidable* — whether a judgement is derivable can be definitively determined by a mechanical process. This is essential for practical implementations of type theories as programming languages, where a type-checking algorithm needs to be able to decide whether a given program is valid.

## 2.3 Type Theory as a Logic

The Curry-Howard correspondence is the observation of the intersection between formal systems of computation and logic. Indeed, the computational calculus of the previous section was defined using inference rules in the style of natural deduction, a proof calculus. More than just an analogy, the correspondence implies that certain computational systems may themselves serve as proof systems.

Our simply-typed  $\lambda$ -calculus coincides with a system of intuitionistic propositional logic. Logical formulae are interpreted as types, with their proofs being the terms of the corresponding types. When a type is inhabited by terms, the corresponding formula may be seen as true, or provable, while uninhabited types correspond to false, or unprovable, formulae. The following table summarises the type-theoretic interpretations for each of the logical constants and connectives of propositional logic.

	Logical Syntax	Type Theory Syntax
False Constant	$\perp$	0
True Constant	Т	1
Disjunction	$P \lor Q$	P+Q
Conjunction	$P \wedge Q$	$P \times Q$
Implication	$P \Rightarrow Q$	$P \to Q$
Negation	$\neg P$	$P\to \mathbb{0}$
Equivalence	$P \Leftrightarrow Q$	$(P \to Q) \times (Q \to P)$

The type  $\mathbb{O}$  is uninhabited, by definition, corresponding to the fact that no proofs of  $\perp$  exist. The  $\mathbb{O}$ -eliminator doubles as the logical principle of explosion, "*ex falso quodlibet*". Conversely, 1 has the constructor  $\star$ , valid in any context, indicating that a proof of  $\top$  may be trivially given at any time. The 1-eliminator corresponds to the fact that assuming  $\top$  as a premise gives a logical argument no additional strength.

Disjunctions are interpreted as sum types — to prove  $P \lor Q$ , it suffices to prove either P, using the inl constructor, or Q, using inr. For a conjunction  $P \land Q$ , proofs for both P and Q must be given, collected into a pair. Implications  $P \Rightarrow Q$  are functions, transforming a proof of P into a proof of Q — the function elimination rule is the logical *modus ponens* rule.

Let us demonstrate some of the laws of propositional logic and their proofs in type theory. We assume a context with arbitrary types P: Type, Q: Type, and R: Type for the rest of this section.

First, some basic properties of implication, the *law of identity*,  $P \Rightarrow P$ , states that implication is a reflexive relation. In addition, implication transitivity is the fact that implications  $P \Rightarrow Q$  and  $Q \Rightarrow R$  may be chained to conclude  $P \Rightarrow R$ .

**Theorem 2.3.1.** Reflexivity,  $P \rightarrow P$ .

Proof.  $\lambda(p:P)$ . p

We are proving an implication by defining a function — we use a  $\lambda$ -term to bind the premise of the implication to a variable, p : P, assuming it as a hypothetical. As the conclusion of the implication is also P, we can simply return the assumption p as the result of the function. The term is both a proof that implication is reflexive, but also the identity function specialised to type P.  $\Box$ 

**Theorem 2.3.2.** Transitivity,  $(P \rightarrow Q) \rightarrow (Q \rightarrow R) \rightarrow P \rightarrow R$ .

Proof.  $\lambda(f: P \to Q)$ .  $\lambda(g: Q \to R)$ .  $\lambda(p: P)$ . g(f p)

Assume hypotheses  $f : P \to Q$ ,  $g : Q \to R$ , and p : P, applying the two functions in turn to conclude g (f p) : R. This proof term corresponds to a higher-order function which composes two argument functions,  $g \circ f$ .

Next, the *law of non-contradiction* asserts a formula cannot be both true and false simultaneously,  $\neg(P \land \neg P)$ , a principle essential for any consistent proof system.

**Theorem 2.3.3.** Non-contradiction,  $(P \times (P \rightarrow \mathbb{Q})) \rightarrow \mathbb{Q}$ .

*Proof.*  $\mathsf{elim}_{\times}$  ( $\lambda(p:P)$ .  $\lambda(np:P \to \mathbb{O})$ .  $np \ p$ )

We are defining a function whose domain is a product,  $P \times (P \to 0)$ , so we apply  $\mathsf{elim}_{\times}$  to decompose it, providing a function to handle the two components of the pair. We bind the those components using  $\lambda$ -terms, so now the context includes p: P and  $np: P \to 0$  with the goal of concluding 0. This is achieved by applying the function np to p.

DeMorgan's laws are good examples which use both product and sum types to represent conjunctions and disjunctions, respectively. **Theorem 2.3.4.** DeMorgan Law,  $(P + Q \rightarrow \mathbb{O}) \rightarrow (P \rightarrow \mathbb{O}) \times (Q \rightarrow \mathbb{O})$ .

 $\textit{Proof. } \lambda(h:P+Q \to \mathbb{O}). \ \langle \lambda(p:P). \ h \ (\mathsf{inl} \ p), \ \lambda(q:Q). \ h \ (\mathsf{inr} \ q) \rangle$ 

Assume  $h: P+Q \to 0$ , and set out to prove both  $P \to 0$  and  $Q \to 0$ , collecting the proofs into a pair. In each case, we bind the hypothesis to a variable with a  $\lambda$ -term and are then to prove 0. This is the conclusion of h, so we can apply it to a value of type P + Q to obtain the desired result — we have p: P in the first case and q: Q in the second, so we use terms inl p and inr q, respectively.  $\Box$ 

Two of the remaining three DeMorgan laws are provable. We give their proof terms here for completeness, but without full explanation.

**Theorem 2.3.5.** DeMorgan Law,  $(P \to \mathbb{O}) \times (Q \to \mathbb{O}) \to P + Q \to \mathbb{O}$ .

*Proof.*  $\mathsf{elim}_{\times}$  ( $\lambda(np: P \to \mathbb{O})$ ).  $\lambda(nq: Q \to \mathbb{O})$ .  $\mathsf{elim}_{+}$  np nq)

**Theorem 2.3.6.** DeMorgan Law,  $(P \to \mathbb{O}) + (Q \to \mathbb{O}) \to P \times Q \to \mathbb{O}$ .

Proof. elim<sub>+</sub> 
$$(\lambda(np: P \to \mathbb{O}). \text{ elim}_{\times} (\lambda(p: P). \lambda(q: Q). np p))$$
  
 $(\lambda(nq: Q \to \mathbb{O}). \text{ elim}_{\times} (\lambda(p: P). \lambda(q: Q). nq q))$ 

However, the final DeMorgan law,  $(P \times Q \to \mathbb{O}) \to (P \to \mathbb{O}) + (Q \to \mathbb{O})$ , cannot be proven in our system. The reason is that the system defined so far represents an intuitionistic propositional logic whose proofs are constructive. When we come to prove the conclusion  $(P \to \mathbb{O}) + (Q \to \mathbb{O})$ , we must choose to either prove the left side with inl or the right side with inr. However, the types P and Qare arbitrary variables which can represent any concrete type — we don't have sufficient information about them to be able to choose a side.

Similarly, the classical law of excluded middle,  $P \vee \neg P$ , and the equivalent double-negation elimination principle,  $\neg \neg P \Rightarrow P$ , are not intuitionistically valid. If excluded middle held in general, the properties of strong normalisation and canonicity guarantee that such a proof would ultimately reduce to inl carrying a proof of P, or inr carrying a proof of  $\neg P$ . In cases where  $P + (P \rightarrow \mathbb{Q})$  can be proven for some type P, we call P decidable, however it is well-known due to Gödel's incompleteness theorems that there will be undecidable propositions in any sufficiently-powerful proof system. If classical reasoning is desired, excluded middle may be postulated by including a term exMid :  $P + (P \rightarrow 1)$  in the context without any defining equations, but the system will lose canonicity as exMid cannot be reduced into its constructors.

Strong normalisation is also important to ensure logical consistency. For example, if we allow general recursive definitions as exist in many practical functional programming languages, we lose strong normalisation. This would enable

construction of a looping term to inhabit any type,  $\Omega = \Omega : \mathbb{O}$ , thereby proving any statement. In Section 2.6, we shall introduce a restricted form of recursion which preserves the wellfoundedness of the reduction relation.

The decidability of the typing and equality judgements allow type theories which correspond to logical systems to be implemented programmatically. This gives rise to type-theoretic *proof assistants* — software systems which provide a language for stating and proving theorems, as well as guiding the user while constructing proofs by showing the current proof goal, the assumptions available in the context, and so on. When the type-checking apparatus of the software decides the constructions of the program to be well-typed, this corresponds to the proofs being checked to be valid.

Some proof assistants, such as Coq [7] and Lean [22], provide a language of tactics, whose terms resemble intuitive proof steps, for example, "assume the premise as hypothesis h" and "split h into its possible cases". Internally, these tactics are transformed into terms of the  $\lambda$ -calculus they are based on, and then type checking proceeds as normal. In other proof assistants, such as Agda [48], proofs are written directly in the primary language, but they implement language servers to provide interactive feedback during proof construction. Idris [9] is another implementation of type theory which may function as a proof assistant, but is more suited for leveraging the type system in practical programming tasks rather than proving mathematical theorems. Many of the results in this thesis have been formalised in Agda, with source code listings available in Appendix A.

## 2.4 Dependent Types

The simply-typed  $\lambda$ -calculus presented up to now isn't truly a Martin-Löf type theory as it lacks *dependent types*, allowing types to depend on values. Types aren't currently first-class entities, the rules defining the system enforcing a strict demarcation between types and their values — they cannot be stored in pairs, or returned from functions, and so on. This is because the type universe Type hasn't itself been given a type.

It is tempting give a rule stating Type : Type, but this leads to a logical inconsistency known as Girard's paradox [30, 36], analogous to Russell's paradox in set theory. To solve this, the idea of stratification is borrowed from Russell's theory of types [44] — there is an infinite, enumerable hierarchy of type universes, each assigned the type of the next.

$$\mathsf{Type}_0 : \mathsf{Type}_1 : \mathsf{Type}_2 : \mathsf{Type}_3 : \mathsf{Type}_4 : \dots$$

When multiple universe levels are used in the formation of a type, the result is at the universe level that is their least upper bound. We won't always be explicit about the levels of type universes in future constructions. That is, T: Type, without specifying an index, is taken to be *universe polymorphic*, valid at any universe level. Most often in this thesis, the lowest universe Type<sub>0</sub> will suffice, but it will be highlighted when a higher level is required.

A function whose codomain is a type universe is called a *type family*. Given A: Type and  $B: A \to \mathsf{Type}$ , the result when applying type family B to a term a: A is a type which may depend on the precise value of a. For example, we can create a type family isTrue :  $2 \to \mathsf{Type}$  which is inhabited when applied to true, but uninhabited for false.

$$\mathsf{isTrue} = \mathsf{elim}_+ \; (\mathsf{elim}_1 \; 1) \; (\mathsf{elim}_1 \; \mathbb{0}) : 2 o \mathsf{Type}$$

The logical interpretation of type families is that they correspond to predicates and relations. In the above case, isTrue true reduces to 1 and so may be trivially proven by  $\star$ , while isTrue false reduces to 0 and cannot be proven in an empty context.

We now introduce  $\Pi$ -types and  $\Sigma$ -types, utilising type families to define typelevel connectives which bind terms. The  $\Pi$ -type can be seen as a dependent version of the function type, where the type of the output may depend on the value of the input. The following rules subsume those given for non-dependent functions.

$$\Pi\text{-}\text{FORM} \frac{\Gamma \vdash A : \text{Type} \qquad \Gamma, a : A \vdash B \ a : \text{Type}}{\Gamma \vdash \Pi_{(a:A)} \ B \ a : \text{Type}}$$
$$\Pi\text{-}\text{INTRO} \frac{\Gamma, a : A \vdash b : B \ a}{\Gamma \vdash \lambda(a : A). \ b : \Pi_{(a:A)} \ B \ a}$$
$$\Pi\text{-}\text{ELIM} \frac{\Gamma \vdash f : \Pi_{(a:A)} \ B \ a}{\Gamma \vdash f \ a : B \ a}$$
$$\beta\text{-}\text{REDUCE} \frac{\Gamma \vdash (\lambda(v : A). \ b) : \Pi_{(a:A)} \ B \ a}{\Gamma \vdash (\lambda(v : A). \ b) \ a = b[v/a] : B \ a}$$
$$\eta\text{-}\text{REDUCE} \frac{\Gamma \vdash (\lambda(v : A). \ b) \ a = b[v/a] : B \ a}{\Gamma \vdash (\lambda(v : A). \ f \ v) = f : \Pi_{(a:A)} \ B \ a}$$

The non-dependent function type can be recovered using a constant type family which does not depend on the input value. We may denote this by not using the variable bound by the  $\Pi$ -type former,  $A \to B = \prod_{(a:A)} B$ : Type. Intuitively, the type  $\Pi_{(a:A)} B a$  represents the product of B indexed by all terms a: A. For example,  $\Pi_{(b:2)} B b$  is equivalent to B true  $\times B$  false. In the Curry-Howard correspondence, this is the interpretation of set-theoretic universal quantification,  $\forall_{(a\in A)} B(a)$  — for any term of a: A, the function will provide a corresponding proof of B a specialised to the input value.

One practical use of  $\Pi$ -types is defining *parametrically polymorphic* functions, that is, functions which operate uniformly on any type. For example, we've previously seen identity functions specialised to individual types, such as  $id_2$ . Similarly, we defined projection functions  $\pi_1$  and  $\pi_2$  for product types, however a context including types A and B was assumed. We can now define these functions universally, in an empty context — the argument types are parametrised by dependent functions.

$$\begin{split} &\mathsf{id} = \lambda(A:\mathsf{Type}).\ \lambda a.\ a &: \Pi_{(A:\mathsf{Type})}\ A \to A \\ &\pi_1 = \lambda(A:\mathsf{Type}).\ \lambda(B:\mathsf{Type}).\ \mathsf{elim}_\times\ (\lambda a.\ \lambda b.\ a): \Pi_{(A:\mathsf{Type})}\ \Pi_{(B:\mathsf{Type})}\ A \times B \to A \\ &\pi_2 = \lambda(A:\mathsf{Type}).\ \lambda(B:\mathsf{Type}).\ \mathsf{elim}_\times\ (\lambda a.\ \lambda b.\ b): \Pi_{(A:\mathsf{Type})}\ \Pi_{(B:\mathsf{Type})}\ A \times B \to B \end{split}$$

It is common to leave some of the parameters implicit in cases where they provide no extra clarity. In addition, we may place a function's parameters at the left side of a definition with the function being applied to them, writing  $f = \lambda v$ . b as f v = b. For example, function composition can be defined generically without explicitly writing the  $\Pi$ -types or  $\lambda$ -terms that bind A: Type, B: Type, and C: Type.

$$f \circ g = \lambda a. f (g a) : (B \to C) \to (A \to B) \to A \to C$$

The  $\Sigma$ -type is a dependent version of the product type, where the type of second element of the pair is allowed to depend on the value of the first. These rules subsume those previously given for products.

$$\begin{split} \Sigma\text{-}\mathrm{Form}\, \frac{\Gamma\vdash A: \mathsf{Type} \qquad \Gamma, a: A\vdash B: A \to \mathsf{Type}}{\Gamma\vdash \Sigma_{(a:A)} \ B \ a: \mathsf{Type}} \\ \Sigma\text{-}\mathrm{Intro}\, \frac{\Gamma\vdash a: A \qquad \Gamma\vdash b: B \ a}{\Gamma\vdash \langle a, \ b \rangle: \Sigma_{(a:A)} \ B \ a} \end{split}$$

Instead of the  $\Sigma$ -eliminator describing how to compute a value of some arbitrary T: Type, we allow the resulting type to depend on the value of the pair using a type family  $T : (\Sigma_{(a:A)} B a) \to \text{Type}$ . The function f provided to operate on the pair's elements therefore needs to have a  $\Pi$ -type.

$$\Sigma\text{-}\mathrm{ELIM} \; \frac{\Gamma \vdash f: \Pi_{(a:A)} \; \Pi_{(b:B \; a)} \; T \; \langle a, \, b \rangle \qquad \Gamma \vdash p: \Sigma_{(a:A)} \; B \; a}{\Gamma \vdash \mathsf{elim}_{\Sigma} \; f \; p: T \; p}$$

$$\Sigma\text{-REDUCE} \frac{\Gamma \vdash f : \Pi_{(a:A)} \Pi_{(b:B \ a)} T \langle a, b \rangle \qquad \Gamma \vdash a : A \qquad \Gamma \vdash b : B \ a}{\Gamma \vdash \mathsf{elim}_{\Sigma} f \langle a, b \rangle = f \ a \ b : T \langle a, b \rangle}$$

As with  $\Pi$ -types, a non-dependent product type may be recovered using a constant type family,  $A \times B = \Sigma_{(a:A)} B$ : Type. As  $\Pi$ -types represent indexed products,  $\Sigma$ -types represent indexed sums. This is the interpretation of set-theoretic existential quantification,  $\exists_{(a \in A)} B(a)$ , though note that these existence proofs must be constructive.

Due to canonicity, a proof of  $\Sigma_{(a:A)} B a$  reduces to a pair  $\langle a, b \rangle$  where a : A is a witness to the existential. For example,  $\langle \text{true}, \star \rangle$  is a proof for  $\Sigma_{(b:2)}$  isTrue b, with true witnessing the proof of  $\star$  : isTrue true. This rules out existence proofs by contradiction as they are an instance of the classical double-negation elimination principle,  $\neg \neg (\exists_{a \in A} B(a)) \Rightarrow \exists_{a \in A} B(a)$ , which doesn't hold constructively.

There can be some confusion caused by the name "dependent product". It is most often used to refer to the indexed products that are  $\Pi$ -types, but sometimes mistaken to mean the dependent generalisation of product types that are  $\Sigma$ -types. Throughout this thesis, this label isn't used, instead opting for "dependent function" for  $\Pi$ -types and "dependent pair" for  $\Sigma$ -types. The type of non-dependent pairs will remain the "product type".

We replace the non-dependent elimination and reduction rules of the other basic types with dependent versions which output in a type family indexed by their terms. As with the other rule upgrades, the originals can be reconstructed by instantiating T with constant type families.

$$\begin{split} \mathbb{Q}\text{-ELIM} & \frac{\Gamma \vdash e: \mathbb{Q}}{\Gamma \vdash \text{elim}_{0} \ e: T \ e} \\ \mathbb{I}\text{-ELIM} & \frac{\Gamma \vdash t: T \star}{\Gamma \vdash \text{elim}_{1} \ t \ u: T \ u} \\ \mathbb{I}\text{-REDUCE} & \frac{\Gamma \vdash t: T \star}{\Gamma \vdash \text{elim}_{1} \ t \ u: T \ u} \\ \mathbb{I}\text{-REDUCE} & \frac{\Gamma \vdash t: T \star}{\Gamma \vdash \text{elim}_{1} \ t \ \star = t: T \ \star} \\ & + \text{-ELIM} & \frac{\Gamma \vdash g: \Pi_{(a:A)} \ T \ (\text{inl } a)}{\Gamma \vdash \text{elim}_{+} f \ g \ s: T \ \sigma} \\ & + \text{-REDUCE}_{\text{inl}} & \frac{\Gamma \vdash f: \Pi_{(c:A)} \ T \ (\text{inl } c)}{\Gamma \vdash \text{elim}_{+} f \ g \ (\text{inl } b)} \\ & + \text{-REDUCE}_{\text{inl}} & \frac{\Gamma \vdash g: \Pi_{(b:B)} \ T \ (\text{inl } b)}{\Gamma \vdash \text{elim}_{+} f \ g \ (\text{inl } a)} \\ & + \text{-REDUCE}_{\text{inl}} & \frac{\Gamma \vdash f: \Pi_{(a:A)} \ T \ (\text{inl } a)}{\Gamma \vdash \text{elim}_{+} \ f \ g \ (\text{inl } a)} \\ & + \text{-REDUCE}_{\text{inl}} & \frac{\Gamma \vdash f: \Pi_{(a:A)} \ T \ (\text{inl } a)}{\Gamma \vdash \text{elim}_{+} \ f \ g \ (\text{inl } a)} \\ & = f \ a: T \ (\text{inl } a) \\ & + \text{-REDUCE}_{\text{inr}} & \frac{\Gamma \vdash g: \Pi_{(a:A)} \ T \ (\text{inl } a)}{\Gamma \vdash \text{elim}_{+} \ f \ g \ (\text{inr } b)} = g \ b: T \ (\text{inr } b)} \end{split}$$

The dependent eliminators are more powerful reasoning principles than their non-dependent counterparts. Rather than just allowing for computation with the types, these eliminators facilitate proving predicates over them by case analysis. Like the non-dependent case, the  $\mathbb{O}$ -eliminator proves any proposition from an inconsistent assumption. To prove a property over 1, it suffices to consider the  $\star$  case, while a sum has the cases of inl and inr to consider. The reduction rules are as before, with updated types.

With dependent types providing connectives for universal and existential quantification, our system now corresponds to a system of intuitionistic higherorder logic. Let us review how the concepts introduced align with those of a set theory with an associated predicate calculus, such as the widely-used Zermelo-Fraenkel system. While there are many points of analogy between the two foundations, there are also significant differences.

In set theory, objects are collected into sets, while type theory collects objects into types. An object in set theory may simultaneously be an element of many sets, with set membership  $a \in A$  being a provable proposition within the logic. Each object in type theory, however, is only assigned to a unique type, known from the moment of the object's creation, a : A being a judgement in the metatheory with no way to express this as a proposition within the type theory itself.

In set theory, the only fundamental construction is the set, with everything else encoded by combining sets. For example, functions are often encoded as sets of pairs which relate the input of the function to the output. In type theory, while encodings may be used, it is also possible to define new type primitives which more clearly express their intent. For example, the natural numbers may be encoded as Church numerals,  $\mathbb{N} = \prod_{(A:\mathsf{Type})} (A \to A) \to A \to A : \mathsf{Type}$ , using the dependent function primitive, but we will see their direct definition using an inductive type in Section 2.6.

Typically, set-theoretic logics are classical, admitting the excluded middle and double-negation elimination reasoning principles, allowing proofs by contradiction. Type theory, however, is a computational calculus as well as a logic, with canonicity providing a strong guarantee about the result of reducing a term. This makes the logic inherently constructive, where disjunction proofs must reduce a proof for one of the two sides, and existence proofs must be accompanied by an algorithm for constructing a witness.

In set theory, there is a separation between the language of the logic and the language of the objects the logic describes. In type theory, propositions *are* types, with no distinction made between types inhabited by data structures and types inhabited by proofs — data and proofs are the same class of object, built from the same syntax. A term of type  $\Pi_{(n:\mathbb{N})} \Sigma_{(p:\mathbb{N})}$  is Prime  $p \times (p > n)$  is *both* a proof that there are infinitely many primes as well as as a function that computes a prime larger than any given input.

## 2.5 Equality and Equivalence

There is currently no way to state or prove the proposition that two terms are equal within the system of type theory we have defined. The only notion of equality seen so far is judgemental, part of the external metatheory used to define the rules of the type theory. We introduce *propositional equality*, also called the *identity type*, as a binary relation internal to the type theory.

$$\equiv \operatorname{-Form} \frac{\Gamma \vdash a : A \quad \Gamma \vdash b : A}{\Gamma \vdash a \equiv b : \mathsf{Type}} \equiv \operatorname{-Intro} \frac{\Gamma \vdash a : A}{\Gamma \vdash \mathsf{refl}_{\equiv} a : a \equiv a}$$

Here,  $a \equiv b$  is the type of proofs for the proposition that terms a and b are equal. The type has just one constructor,  $\operatorname{refl}_{\equiv}$ , which proves the reflexive case  $a \equiv a$ .

Judgemental equality implies propositional equality — judgementally-equal terms can always be reduced to their identical canonical forms, the equality then witnessed by  $refl_{\equiv}$ . As an example illustrating this, we can show that function composition is associative.

**Theorem 2.5.1.** Associativity,  $(f \circ (g \circ h)) \ a \equiv ((f \circ g) \circ h) \ a$ .

*Proof.*  $\operatorname{assoc}_{\circ} = \operatorname{refl}_{\equiv} (f (g (h a)))$ 

By applying the definition of function composition, both sides of the equation reduce to f(g(h a)) by judgemental equalities, so it can be proved directly by the refl<sub>=</sub> constructor.

Propositional equality is more powerful than judgemental equality, allowing us to prove terms equal which are not necessarily judgementally equal, however propositional equality is not decidable in general. We can see an application of the dependent 1-eliminator in combination with propositional equality by proving that  $\star$  is the only canonical term for 1. Similarly, every Boolean value is either equal to **true** or **false**, whose proof also uses the dependent eliminator for sums.

**Theorem 2.5.2.** Uniqueness,  $\Pi_{(u:1)} u \equiv \star$ .

*Proof.*  $\mathsf{elim}_1$  ( $\mathsf{refl}_{\pm} \star$ )

Applying the eliminator, we are only required to prove the case where u is  $\star$ . The conclusion becomes  $\star \equiv \star$ , proved by  $\mathsf{refl}_{\equiv} \star$ .

Lemma 2.5.3.  $\Pi_{(b:2)}$   $(b \equiv true) + (b \equiv false)$ 

*Proof.*  $\mathsf{elim}_+$  ( $\mathsf{elim}_1$  (inl ( $\mathsf{refl}_{\equiv}$  true))) ( $\mathsf{elim}_1$  (inr ( $\mathsf{refl}_{\equiv}$  false)))

Since 2 is defined as 1+1, the sum-eliminator splits the proof into the inl case and the inr case. In either case, the 1-eliminator has us consider  $\star$  the only possible term being carried by each constructor. Considering the judgemental equality defining true = inl  $\star$ , in the first case the conclusion is (true  $\equiv$  true)+(true  $\equiv$  false), so we decide to prove the left side of the sum with inl (refl $\equiv$  true). In the second case, we are to prove (false  $\equiv$  true) + (false  $\equiv$  false), so we elect for the right side this time.

Lemma 2.5.3 validates there are no Boolean terms beside those equal to true or false, but it doesn't establish that these are actually distinct. For that, we need to prove true  $\equiv$  false  $\rightarrow 0$ , that true and false being equal leads to absurdity.

To prove statements with a propositional equality as a hypothesis, we need an elimination rule. Since the equality type former is a binary type family, the output predicate needs to be a dependent function which takes the parameters of the family,  $T: \prod_{(a,b:A)} a \equiv b \rightarrow \mathsf{Type}.$ 

$$\equiv \text{-ELIM } (J) \frac{\Gamma \vdash f : \Pi_{(c:A)} T c c (\operatorname{refl}_{\equiv} c) \qquad \Gamma \vdash p : a \equiv b}{\Gamma \vdash \operatorname{elim}_{\equiv} f a b p : T a b p}$$

$$\equiv -\text{REDUCE} \frac{\Gamma \vdash f : \Pi_{(c:A)} \ T \ c \ c \ (\text{refl}_{\equiv} \ c)}{\Gamma \vdash \text{elim}_{\equiv} \ f \ a \ a \ (\text{refl}_{\equiv} \ a) = f \ a : T \ a \ a \ (\text{refl}_{\equiv} \ a)}$$

The elimination rule, known as the *J*-rule, states that to prove a predicate on a equality proof  $a \equiv b$ , it suffices to prove the predicate in the reflexive case  $\operatorname{refl}_{\equiv} c : c \equiv c$ , for an arbitrary c of the type.

We need some basic properties of equality before we can prove that true and false are distinct, however. Let us start with symmetry, which instantiates the *J*-rule with predicate  $T \ a \ b \ p = b \equiv a$ .

**Theorem 2.5.4.** Symmetry,  $a \equiv b \rightarrow b \equiv a$ .

*Proof.* sym<sub> $\equiv$ </sub> = elim<sub> $\equiv$ </sub> ( $\lambda(c:A)$ . refl<sub> $\equiv$ </sub> c)

We assume  $a \equiv b$  as a hypothesis, and are to prove  $b \equiv a$ . Applying the eliminator to the hypothesis, we need to provide a function of variable c which proves  $T \ c \ c \ (refl_{\equiv} c)$ . Reducing this expression using the definition of T, this is simply the statement  $c \equiv c$ , which can be easily discharged using the  $refl_{\equiv}$  constructor. In essence, since we started from an assumption of  $a \equiv b$ , the J-rule allowed us to substitute all of the instances of a and b by a new variable c in the conclusion we wanted to reach.

Similarly, transitivity can be proved by fixing c : A and taking the predicate to be  $T \ a \ b \ p = b \equiv c \rightarrow a \equiv c$ .

**Theorem 2.5.5.** Transitivity,  $a \equiv b \rightarrow b \equiv c \rightarrow a \equiv c$ .

*Proof.* trans<sub> $\equiv$ </sub> = elim<sub> $\equiv$ </sub> ( $\lambda(d:A)$ . id<sub>( $d\equiv c$ )</sub>)

Assume  $a \equiv b$  as a hypothesis, but leave  $b \equiv c$  unassumed — we must prove the implication  $b \equiv c \rightarrow a \equiv c$ . Apply the eliminator and instantiate the predicate with new variable d, leaving us to prove  $d \equiv c \rightarrow d \equiv c$ . This can be proven with an identity function id specialised to type  $d \equiv c$ 

A common pattern in functional programming, and in mathematics more broadly, is theorem proving by *equational reasoning*, for example, sequencing equality proofs  $a \equiv b \equiv \ldots \equiv y \equiv z$  to prove  $a \equiv z$ . This is justified in our type theory by chaining applications of trans<sub> $\equiv$ </sub>.

As a final basic property we would expect an equality to exhibit, propositional equality forms a congruence relation. Sometimes called a substitution or transportation principle, this ensures that if we know  $a \equiv b$ , we may substitute b for a while preserving the proof of any predicate,  $P \ a \rightarrow P \ b$ . Taken in combination with symmetry, the direction of the implication is unimportant and the two terms become interchangeable.

**Theorem 2.5.6.** Congruence,  $P \ a \rightarrow a \equiv b \rightarrow P \ b$ .

*Proof.* cong  $(P : A \to \mathsf{Type}) = \lambda(p : P a)$ .  $\lambda(q : a \equiv b)$ .  $\mathsf{elim}_{\equiv} (\lambda(c : A), \mathsf{id}_{(P c)}) q p$ 

Resembling the proof of transitivity, after applying the eliminator the conclusion becomes  $P \ c \to P \ c$ , proved by an identity function.

A useful special case of congruence is that functions respect equality. That is, when proving an equality where both sides are applications of the same function, the function may be stripped away, simplifying the equation to be proved.

**Theorem 2.5.7.** Functions respect equality,  $a \equiv b \rightarrow f \ a \equiv f \ b$ .

*Proof.* resp  $(f : A \to B) = \text{elim}_{\equiv} (\lambda(c : A), \text{ refl}_{\equiv} (f c))$ 

Using the congruence property with the previously-defined isTrue predicate which maps true  $\mapsto 1$  and false  $\mapsto 0$ , we can prove that 2 is truly a Boolean type.

#### Lemma 2.5.8. true $\equiv$ false $\rightarrow 0$

#### *Proof.* cong isTrue $\star$

Assume as a hypothesis that true  $\equiv$  false and set out to prove  $\mathbb{O}$ , which is judgementally equal to isTrue false. Since we have true  $\equiv$  false by assumption, congruence tells us that we may substitute within the conclusion, instead proving isTrue true, which is judgementally equal to 1. This is achieved with the trivial proof  $\star$ .

**Theorem 2.5.9.** Boolean,  $(\Pi_{(b:2)} \ (b \equiv \mathsf{true}) + (b \equiv \mathsf{false})) \times (\mathsf{true} \equiv \mathsf{false} \to \mathbb{O}).$ 

*Proof.* This is a straightforward pairing of the proof terms for Lemmas 2.5.3 and 2.5.8.  $\hfill \square$ 

As types are first-class values, it is also natural to ask when two types are equal. For types A and B, the type  $A \equiv B$  contains proofs of their identification. Of course, judgementally-equal types can be proved propositionally equal by reflexivity, for example,  $\operatorname{refl}_{\equiv} 2 : 2 \equiv 2$ . However, in intensional type theories, such as the one defined in this chapter, we cannot go beyond judgemental equality of types and must instead talk of their equivalence.

The idea of equivalence has been used informally up until now, actually referring to three distinct notions. The first is propositional equivalence, corresponding to logical equivalence,  $A \leftrightarrow B = (A \rightarrow B) \times (B \rightarrow A)$ , indicating that the types are coinhabited — either they are both empty or they each have at least one term. However, this says nothing about the number of distinct terms inhabiting non-empty types, for example,  $1 \leftrightarrow 2$  is itself inhabited.

Second, it has previously been noted that some types, such as 1 + (1 + 1) and (1 + 1) + 1, are equivalent in the sense that they have the same number of canonical terms despite being structurally distinct. Let us now be more precise about what this form of type equivalence entails.

**Definition 2.5.10.** Two types A: Type and B: Type are equivalent when there is a bijective mapping between their terms — there are functions  $f : A \to B$  and  $g : B \to A$  which are inverse.

$$A \simeq B = \Sigma_{(f:A \to B)} \Sigma_{(g:B \to A)} (\Pi_{(a:A)} \ a \equiv (g \circ f) \ a))$$
$$\times (\Pi_{(b:B)} \ b \equiv (f \circ g) \ b))$$

A proof of type equivalence is therefore a term whose canonical form consists of nested pairs,  $\langle f, \langle g, \langle p, q \rangle \rangle$ , containing two functions and two equality proofs. For notational convenience, we shall write this as a quadruple,  $\langle f, g, p, q \rangle$ .

While we cannot prove 1 + (1 + 1) and (1 + 1) + 1 equal in an intensional type theory, we can prove them equivalent.

Theorem 2.5.11.  $1 + (1 + 1) \simeq (1 + 1) + 1$ 

#### Proof. $\langle f, g, p, q \rangle$

Where f maps inl  $\star \mapsto \text{inr } \star$ , inr (inl  $\star$ )  $\mapsto$  inl (inr  $\star$ ), inr (inr  $\star$ )  $\mapsto$  inl (inl  $\star$ ), and g is its inverse. To be fully-formal, the terms for the function definitions, as well as proof terms p and q, would be be defined by case analysis, using the eliminator  $\text{elim}_+$ . This is one possible equivalence proof out of a total of six for this type. The others may be produced by permuting the terms in the output of the map.  $\Box$ 

Of course, like propositional equality, type equivalence forms an equivalence relation.

#### **Theorem 2.5.12.** Reflexivity, $A \simeq A$ .

*Proof.*  $\operatorname{refl}_{\simeq} = \langle \operatorname{id}_A, \operatorname{id}_A, \operatorname{refl}_{\equiv}, \operatorname{refl}_{\equiv} \rangle$ 

The identity function  $\mathsf{id}_A$  is used to instantiate both functions, so the final two elements must be of type  $\Pi_{(a:A)} a \equiv \mathsf{id}_A (\mathsf{id}_A a)$ . Since  $\mathsf{id}_A (\mathsf{id}_A a)$  reduces to a, this is judgementally equal to the type of the  $\mathsf{refl}_{\equiv}$  constructor.

**Theorem 2.5.13.** Symmetry,  $A \simeq B \rightarrow B \simeq A$ .

*Proof.* sym<sub> $\simeq$ </sub> = elim<sub> $\Sigma$ </sub> ( $\lambda f$ . elim<sub> $\Sigma$ </sub> ( $\lambda g$ . elim<sub> $\Sigma$ </sub> ( $\lambda p$ .  $\lambda q$ . (g, f, q, p))))

Assume the equivalence  $\langle f, g, p, q \rangle$  and simply swap the functions and swap the proofs, resulting in  $\langle g, f, q, p \rangle$ .

#### **Theorem 2.5.14.** Transitivity, $A \simeq B \rightarrow B \simeq C \rightarrow A \simeq C$ .

*Proof.* trans<sub> $\simeq$ </sub>  $(d: A \simeq B)$   $(e: B \simeq C) = \dots$ 

Assume the two equivalences  $\langle f, g, p, q \rangle$  and  $\langle f', g', p', q' \rangle$ . To construct the two functions for the output equivalence, we compose the corresponding functions from the hypotheses,  $f' \circ f$  and  $g \circ g'$ . For the proofs, we need to show, for arbitrary a and c, that  $a \equiv ((g \circ g') \circ (f' \circ f)) a$  and  $c \equiv ((f' \circ f) \circ (g \circ g')) c$ .

With the associative property, we can group the compositions such that our equality proofs p, q, p', and q' can eliminate them.

		$((g \circ g') \circ (f' \circ f)) a$			$((f' \circ f) \circ (g \circ g')) \ c$
$(assoc_\circ)$	≡	$(g \circ (g' \circ f') \circ f) \ a$	$(assoc_\circ)$	≡	$(f' \circ (f \circ g) \circ g') \ c$
(p')	≡	$(g \circ f) a$	(q)	≡	$(f' \circ g') \ c$
(p)	≡	a	(q')	≡	c

The full proof term is omitted here to avoid the tedious decomposition of the two equivalences into their eight components using six applications of  $\operatorname{elim}_{\Sigma}$ . Proof assistants provide *pattern matching* for cases like this, as seen in the version of this proof included in Appendix A.

The third notion of equivalence is equivalence of functions. For example, bubble-sort and merge-sort can be considered equivalent functions in the sense they produce the same outputs when given the same inputs. However, since their implementations differ vastly, they cannot be proven propositionally equal. **Definition 2.5.15.** Two functions  $f : \prod_{(a:A)} B a$  and  $g : \prod_{(a:A)} B a$  are equivalent when they are pointwise equal.

$$f \sim g = \prod_{(a:A)} f \ a \equiv g \ a$$

Naturally, this also forms an equivalence relation. The proofs are immediate from the reflexivity, symmetry, and transitivity of the underlying propositional equality relation.

**Theorem 2.5.16.** Reflexivity,  $f \sim f$ .

*Proof.*  $\operatorname{refl}_{\sim}(a:A) = \operatorname{refl}_{\equiv}(f a)$ **Theorem 2.5.17.** Symmetry,  $f \sim g \rightarrow g \sim f$ . *Proof.*  $\operatorname{sym}_{\sim} (e: f \sim g) (a: A) = \operatorname{sym}_{=} (e \ a)$ 

**Theorem 2.5.18.** Transitivity,  $f \sim g \rightarrow g \sim h \rightarrow f \sim h$ .

Proof. trans<sub>~</sub>  $(d: f \sim g) (e: g \sim h) (a: A) = \text{trans}_{\equiv} (d a) (e a)$ 

As corollaries from their reflexivity, both type equivalence and function equivalence can be shown to be implied by propositional equality. Simply use congruence to substitute one of the sides of the conclusion so it becomes an instance of reflexivity.

**Corollary 2.5.19.** Type equality implies equivalence,  $A \equiv B \rightarrow A \simeq B$ .

*Proof.* cong ( $\lambda(C : \mathsf{Type})$ .  $A \simeq C$ ) (refl<sub> $\simeq$ </sub> A)

**Corollary 2.5.20.** Function equality implies equivalence,  $f \equiv g \rightarrow f \sim g$ .

*Proof.* cong  $(\lambda(h: \Pi_{(a:A)} B a). f \sim h)$  (refl<sub>~</sub> f)

Equality and equivalence are very active topics of research in the field of type theory. It is often desirable to eschew *intentional* type theories, where certain equality proofs may be difficult or impossible, and only focus on the *extensional* properties of the objects in question — if two distinct types or functions behave identically, why shouldn't they be considered equal?

For ease of reasoning about functions, the axiom of function extensionality may be assumed. This states the inverse of the above, that function equivalence implies equality, funExt :  $f \sim g \rightarrow f \equiv g$ , allowing us to ignore their internal details and consider only their external behaviour. The axiom is independent of the intensional type theory we've defined — it cannot be proved directly, but its assumption is consistent with the existing rules [54]. However, as with

adding classical reasoning principles as additional axioms, the funExt term has no reduction rules associated with it and so canonicity is lost.

The development of homotopy type theory [55] takes extensionality to its limit with the univalence axiom, due to Vladimir Voevodsky. Types are interpreted as topological spaces, with proofs of equality being paths between their points, and function equivalences regarded as homotopies deforming one function into the other. Notably, univalence implies, in addition to function extensionality, that type equivalence implies equality,  $A \simeq B \rightarrow A \equiv B$  — equivalent types are truly equal, their implementation details irrelevant. There has been a recent effort to give computational content to univalence, one result being the system of cubical type theory [17], where univalence is a theorem rather than an axiom, and canonicity is preserved [35].

Since extensionality is not the focus of this thesis, only a very basic overview has been given. For the most part, we will operate within an intentional type theory, although function extensionality is assumed for Chapter V.

Finally in this section, we summarise the various equality and equivalence relations of type theory, the syntax used for them throughout this thesis, and compare their relative logical strength.

- JUDGEMENTAL EQUALITY, a = b: a statement in the defining metatheory, used to add new definitions to the type theory and to specify how terms are evaluated to their canonical forms. It is the strongest notion of equality, implying every relation on this list. However, it cannot be expressed as a type in the syntax of the type theory itself, so we cannot reason about judgemental equalities from within the system.
- PROPOSITIONAL EQUALITY,  $a \equiv b$ : the notion of equality internal to the system of type theory, able to be stated as a type and reasoned about from within the system. Propositional equality is weaker than judgemental equality, but still implies all of the equivalences lower down on this list.
- TYPE EQUIVALENCE,  $A \simeq B$ : a bijective mapping between the terms of two types, establishing that the types have the same size. Type equivalence is weaker than propositional equality, but stronger than propositional equivalence.
- PROPOSITIONAL EQUIVALENCE, A ↔ B: a mapping between the terms of two types, implied by type equivalence since there is no requirement for injectivity or surjectivity. This establishes that the types are coinhabited, both empty or each inhabited by at least one term. This is the Curry-Howard interpretation of two propositions being logically equivalent.

• FUNCTION EQUIVALENCE,  $f \sim g$ : a statement that two functions are pointwise equal, with equal inputs to the functions yielding equal outputs. Weaker than propositional equality in intensional type theories, but of the same strength as propositional equality in extensional type theories where function extensionality is adopted.

### 2.6 Inductive Types

Combining sum and product types with the basic types  $\mathbb{O}$  and  $\mathbb{1}$ , we are able to form any finite type, but is also possible to encode more complex structures using functions, as Church did in his  $\lambda$ -calculus [14,15]. For example, the Church numerals are an encoding of the natural numbers as parametrically polymorphic functions,  $\mathbb{N} = \prod_{(A:\mathsf{Type})} (A \to A) \to A \to A$ . The numeral representing natural number *n* takes a function as an argument and iterates it *n* times over a second argument.

$$0 = \lambda A. \ \lambda f. \ \lambda x. \ x$$
  

$$1 = \lambda A. \ \lambda f. \ \lambda x. \ f \ x$$
  

$$2 = \lambda A. \ \lambda f. \ \lambda x. \ f \ (f \ x)$$
  

$$3 = \lambda A. \ \lambda f. \ \lambda x. \ f \ (f \ (f \ x)))$$
  
...

It can be shown by parametricity [52, 57] that the only closed terms of this type are extensionally equal to the Church numerals. However, in a setting without function extensionality, it is more convenient to define the natural numbers, and other tree-like structures, as dedicated types with new syntax. In addition, our strictly-enforced type system prevents the definition of recursive functions using fixed-point combinators, such as Haskell Curry's Y combinator,  $\lambda f. (\lambda x. f (x x)) (\lambda x. f (x x))$  — the self-application x x is invalid as it cannot be assigned a type.

We define Peano-style natural numbers in the same way as our other types, with formation, introduction, elimination, and reduction rules.

$$\mathbb{N}\text{-}\mathrm{FORM} \ \overline{\mathbb{N}:\mathrm{Type}} \qquad \mathbb{N}\text{-}\mathrm{INTRO}_{\mathsf{zero}} \ \overline{\mathsf{zero}:\mathbb{N}} \qquad \mathbb{N}\text{-}\mathrm{INTRO}_{\mathsf{succ}} \ \overline{\frac{\Gamma \vdash n:\mathbb{N}}{\Gamma \vdash \mathsf{succ} n:\mathbb{N}}} \\ \mathbb{N}\text{-}\mathrm{ELIM} \ \overline{\frac{\Gamma \vdash z:T \ \mathsf{zero} \qquad \Gamma \vdash s:\Pi_{(n:\mathbb{N})} \ T \ n \to T \ (\mathsf{succ} \ n) \qquad \Gamma \vdash n:\mathbb{N}}{\Gamma \vdash \mathsf{elim}_{\mathbb{N}} \ z \ s \ n:T \ n}} \\ \mathbb{N}\text{-}\mathrm{Reduce}_{\mathsf{zero}} \ \overline{\frac{\Gamma \vdash z:T \ \mathsf{zero} \qquad \Gamma \vdash s:\Pi_{(n:\mathbb{N})} \ T \ n \to T \ (\mathsf{succ} \ n)}{\Gamma \vdash \mathsf{elim}_{\mathbb{N}} \ z \ s \ \mathsf{zero} = z:T \ \mathsf{zero}}}$$
$$\mathbb{N}\text{-}\operatorname{Reduce}_{\mathsf{succ}} \frac{\Gamma \vdash z : T \text{ zero } \Gamma \vdash s : \Pi_{(n:\mathbb{N})} T n \to T (\operatorname{\mathsf{succ}} n) \quad \Gamma \vdash n : \mathbb{N}}{\Gamma \vdash \operatorname{\mathsf{elim}}_{\mathbb{N}} z \ s (\operatorname{\mathsf{succ}} n) = s \ n (\operatorname{\mathsf{elim}}_{\mathbb{N}} z \ s \ n) : T (\operatorname{\mathsf{succ}} n)}$$

The constructors for the natural numbers, defined by the two introduction rules, are zero and succ. Unlike previous introduction rules, the succ rule is recursive, requiring a natural number n to be constructed before the succ constructor may be applied to yield a new natural number. The canonicity principle specialised to natural numbers states that any closed term of type  $\mathbb{N}$  is either the numeral zero, or the succ constructor applied to a numeral.

$$0 = \text{zero}$$
  

$$1 = \text{succ zero}$$
  

$$2 = \text{succ (succ zero)}$$
  

$$3 = \text{succ (succ (succ zero))}$$

Seen as a computation principle, the eliminator gives a mechanism for defining functions by primitive recursion, saying how to compute from **zero** and how to compute from **succ** n, with the reduction rules specifying the recursion scheme. Seen as a proof principle, the elimination rule allows proofs by induction over natural numbers — to prove a predicate T for all natural numbers, it suffices to prove the base case, T **zero**, and the inductive case,  $\Pi_{(n:\mathbb{N})} T n \to T$  (**succ** n).

To be convinced that this type faithfully represents the natural numbers, we would also expect the **succ** constructor to be injective, which is explicitly stated as one of Peano's axioms. We may prove this as a theorem rather than adding an axiom, requiring the definition of a type family indexed by the natural numbers.

Theorem 2.6.1. succ  $m \equiv \text{succ } n \rightarrow m \equiv n$ 

*Proof.* cong (elim<sub>N</sub>  $\mathbb{O}$  ( $\lambda o. \ \lambda r. \ m \equiv o$ )) (refl<sub> $\equiv$ </sub> m)

Using the natural number eliminator, define the type family  $P : \mathbb{N} \to \mathsf{Type}$ which maps succ  $o \mapsto m \equiv o$ . The type zero maps to is of no consequence and is chosen arbitrarily.

Specialising cong P to the case where both numbers are successors, its type is P (succ m)  $\rightarrow$  succ  $m \equiv$  succ  $n \rightarrow P$  (succ n). Reducing the applications of Pyields  $m \equiv m \rightarrow$  succ  $m \equiv$  succ  $n \rightarrow m \equiv n$ . After giving a proof of  $m \equiv m$ , this is exactly the statement that succ is injective.

A basic operator that doesn't require recursion is the predecessor function, which removes an application of **succ** from positive natural numbers.

$$\mathsf{pred} = \mathsf{elim}_{\mathbb{N}} \mathsf{zero} (\lambda n. \ \lambda r. \ n) : \mathbb{N} \to \mathbb{N}$$

More complex arithmetic operators may be defined by primitive recursion. For example, an intuitive definition of addition of two natural numbers can be defined by pattern matching on the second argument. In the base case, the first argument is returned, while the recursive case moves the application of **succ** away from the second argument.

 $\begin{array}{l} \mathsf{add}:\mathbb{N}\to\mathbb{N}\to\mathbb{N}\\ \mathsf{add}\ m\ \mathsf{zero} &=m\\ \mathsf{add}\ m\ (\mathsf{succ}\ n)=\mathsf{succ}\ (\mathsf{add}\ m\ n) \end{array}$ 

This definition is realised by applying the eliminator non-dependently, instantiated at constant type family  $T \ n = \mathbb{N}$ .

$$\mathsf{add} = \lambda m. \mathsf{elim}_{\mathbb{N}} \ m \ (\lambda n. \ \lambda r. \ \mathsf{succ} \ r) : \mathbb{N} \to \mathbb{N} \to \mathbb{N}$$

The first argument to to  $\operatorname{elim}_{\mathbb{N}}$  represents the base case, while the second argument is the recursive case, with the result of the recursive call being bound to variable r. Applying the reduction rule for succ, this can be seen to reduce to the expected expression.

add 
$$m$$
 (succ  $n$ ) =  $\operatorname{elim}_{\mathbb{N}} m$  ( $\lambda n. \lambda r. \operatorname{succ} r$ ) (succ  $n$ )  
 $\rightarrow$  ( $\lambda n. \lambda r. \operatorname{succ} r$ )  $n$  ( $\operatorname{elim}_{\mathbb{N}} m$  ( $\lambda n. \lambda r. \operatorname{succ} r$ )  $n$ )  
 $\rightarrow^* \operatorname{succ} (\operatorname{elim}_{\mathbb{N}} m (\lambda n. \lambda r. \operatorname{succ} r) n)$   
= succ (add  $m n$ )

As addition was defined by recursion on the second argument, the fact that zero is the right identity for addition is a judgemental equality — the proof of add  $n \text{ zero } \equiv n$  is simply  $\text{refl}_{\equiv} n$ . However, for the left identity, add zero n cannot be reduced without knowing whether n is constructed by zero or succ. To show this, we use the eliminator dependently, instantiating the predicate with  $T n = \text{add zero } n \equiv n$ , as a proof by induction.

**Theorem 2.6.2.**  $\Pi_{(n:\mathbb{N})}$  add zero  $n \equiv n$ 

*Proof.* elim<sub>N</sub> (refl<sub> $\equiv$ </sub> zero) ( $\lambda n$ .  $\lambda r$ . resp succ r)

Proceed by induction on n. The base case is add zero zero  $\equiv$  zero. By definition of add, this reduces to zero  $\equiv$  zero, which is proved by refl $_{\equiv}$  zero. This therefore becomes the first argument to  $\operatorname{elim}_{\mathbb{N}}$ .

In the inductive case, we have inductive hypothesis add zero  $n \equiv n$  and must show add zero (succ n)  $\equiv$  succ n, which reduces to succ (add zero n)  $\equiv$  succ n. Since succ respects equality, by Theorem 2.5.7, it can be stripped from both sides of the conclusion, with the remaining term being identical to the inductive hypothesis. In the proof term, it can be seen that the inductive case is a dependent function where variable r binds the inductive hypothesis.

Generalising from the definition of natural numbers, we can define an entire class of types called *inductive types* whose terms are wellfounded trees. However, instead of explicitly specifying the rules for all these types, we use a more concise notation which specifies the type of the former, and the types of all constructors.

inductive 
$$\mathbb{N}$$
 : Type  
zero :  $\mathbb{N}$   
succ :  $\mathbb{N} \to \mathbb{N}$ 

This definition of natural numbers is understood to be shorthand for all of its rules. That is, a type formation rule, an introduction rule for each constructor, an elimination rule, and a reduction rule for each constructor.

We can define inductive type families of arbitrary arity by setting the type former to be a function type. Each constructor may also be of arbitrary arity, and may depend recursively on other elements of the type, forming new branches in the tree. These types may be given as instances of W-types [45], but we will opt to allow the direct definition of recursive functions by pattern matching [23].

Natural numbers are unary trees whose nodes are applications of the succ constructor and whose leaves are zero. The tree carries no additional data at the nodes aside from subtrees, the only information represented being the tree's depth. If we instead define the inductive type of unary trees where the nodes additionally carry a term of type A, we obtain the List data structure.

```
inductive List : Type \rightarrow Type

[] : \Pi_{(A:\mathsf{Type})} List A

:: : \Pi_{(A:\mathsf{Type})} A \rightarrow \mathsf{List} A \rightarrow \mathsf{List} A
```

List is defined as a type family, giving dependent function types for its constructors to allow for parametric polymorphism, but dependent types may also be used to impose constraints on when the use of a constructor is valid. We can define the Fin inductive family of finite types, indexed by natural numbers, where Fin n has exactly n canonical terms, representing those natural numbers strictly less than n.

 $\begin{array}{l} \text{inductive Fin} : \mathbb{N} \to \mathsf{Type} \\ \\ \text{zero}_{\mathsf{Fin}} : \Pi_{(n:\mathbb{N})} \; \mathsf{Fin} \; (\mathsf{succ} \; n) \\ \\ \\ \text{succ}_{\mathsf{Fin}} : \Pi_{(n:\mathbb{N})} \; \mathsf{Fin} \; n \to \mathsf{Fin} \; (\mathsf{succ} \; n) \end{array}$ 

Since both constructors result in Fin (succ n), there is no way of constructing a term of type Fin 0. The only way of constructing a term of type Fin 1 is using  $zero_{Fin}$ , as  $succ_{Fin}$  would have be applied to a term of type Fin 0. The  $succ_{Fin}$  constructor may be used for Fin 2 onwards, but only a limited number of times before  $zero_{Fin}$  is forced.

Inductive families representing predicates and relations may also be defined. For example, the ordering of natural numbers can be established with an inductive binary relation.

$$\begin{array}{l} \mathsf{inductive} \leq : \mathbb{N} \to \mathbb{N} \to \mathsf{Type} \\ \mathsf{zero}_{\leq} \, : \, \Pi_{(n:\mathbb{N})} \; \mathsf{zero} \leq n \\ \mathsf{succ}_{\leq} \, : \, \Pi_{(m,n:\mathbb{N})} \; m \leq n \to \mathsf{succ} \; m \leq \mathsf{succ} \; n \end{array}$$

The strength of inductive types is that, in lieu of an elimination rule, we allow direct definitions by *structural recursion*, pattern matching on the possible ways of constructing a type. To be certain of normalisation, recursive calls must be on a structurally smaller term than the original argument, that is a subexpression with strictly fewer constructors.

This is more general than the primitive recursion allowed by  $\operatorname{elim}_{\mathbb{N}}$ , allowing multiple constructors to be stripped at once. For example, we can define the function which computes  $n \mapsto \lfloor n/2 \rfloor$  by pattern matching two constructors deep.

$$\begin{split} \mathsf{half} &: \mathbb{N} \to \mathbb{N} \\ \mathsf{half zero} &= \mathsf{zero} \\ \mathsf{half} (\mathsf{succ zero}) &= \mathsf{zero} \\ \mathsf{half} (\mathsf{succ} (\mathsf{succ} n)) &= \mathsf{succ} (\mathsf{half} n) \end{split}$$

Seen as universally-quantified proofs, dependent functions defined by structural recursion correspond to proofs by structural induction. For example, the reflexivity of  $\leq$  can be established by induction using the eliminator  $\mathsf{elim}_{\mathbb{N}}$ , but we demonstrate the proof as a structurally-recursive function.

**Theorem 2.6.3.** Reflexivity,  $n \leq n$ .

Proof.

$$\begin{split} \operatorname{refl}_{\leq} &: \Pi_{(n:\mathbb{N})} \ n \leq n \\ \operatorname{refl}_{\leq} \operatorname{zero} &= \operatorname{zero}_{\leq} \\ \operatorname{refl}_{\leq} (\operatorname{succ} n) = \operatorname{succ}_{\leq} (\operatorname{refl}_{\leq} n) \end{split}$$

Proceed by induction on n. In the base case, we have to show  $zero \leq zero$ , established by the  $zero_{\leq}$  constructor which states that zero is less than or equal any other natural number.

In the inductive case, we have inductive hypothesis  $n \leq n$  and need to prove succ  $n \leq \text{succ } n$ . The succ<sub> $\leq$ </sub> constructor states that the relation respects succ, so we are just left to prove a statement equal to the inductive hypothesis. To access the inductive hypothesis in the proof term we call  $refl_{\leq}$  recursively, at a structurally-smaller index.

In the presence of inductive families with dependently-typed constructors, structural recursion becomes a powerful reasoning tool, allows patterns to be refined according to the type of the constructor [18,31]. We demonstrate this by proving the transitivity of  $\leq$  by performing structural recursion on a proof of the relation rather than a natural number.

**Theorem 2.6.4.** Transitivity,  $m \le n \to n \le o \to m \le o$ .

Proof.

$$\begin{split} & \operatorname{trans}_{\leq}: \Pi_{(m,n,o:\mathbb{N})} \ m \leq n \to n \leq o \to m \leq o \\ & \operatorname{trans}_{\leq} \operatorname{zero}_{\leq} q = \operatorname{zero}_{\leq} \\ & \operatorname{trans}_{\leq} (\operatorname{succ}_{\leq} p) \ (\operatorname{succ}_{\leq} q) = \operatorname{succ}_{\leq} (\operatorname{trans}_{\leq} p \ q) \end{split}$$

Assume  $m \leq n$  and  $n \leq o$ , towards a proof of  $m \leq o$ . Proceed by structural induction on the proof of  $m \leq n$ .

In the base case, it is constructed by  $\text{zero}_{\leq} : \text{zero} \leq n$ , so the only possibility is when *m* is zero. The conclusion is therefore refined to zero  $\leq o$ , which is established by  $\text{zero}_{\leq}$ .

In the inductive case, the proof is constructed by  $\operatorname{succ}_{\leq} p : \operatorname{succ} m \leq \operatorname{succ} n$ , the only possible pattern is when both sides of the relation are successors. The type of our second assumption is refined to  $\operatorname{succ} n \leq o$ , which forces it to be constructed by  $\operatorname{succ}_{\leq} q : \operatorname{succ} n \leq \operatorname{succ} o$ , the third number must also be a successor. In light of these refinements, the conclusion is  $\operatorname{succ} m \leq \operatorname{succ} o$ , so we can apply the  $\operatorname{succ}_{\leq}$  constructor and then prove  $m \leq o$  by recursion on the structurally smaller terms p and q.

# 2.7 Coinductive Types

Inductive types represent wellfounded trees, structural recursion restricting them to finite depth, however infinite structures are often objects of study. In functional programming languages featuring general recursion and lazy evaluation, such as Haskell [42], no distinction is made between inductive types with finite objects and those with infinite objects [34]. For example, a definition like the following would be valid, defining the infinite list nats = 0 :: 1 :: 2 :: ... of all natural numbers — in fact, this is so common in Haskell that the shorthand [0..] is provided

nats : List  $\mathbb{N}$ from :  $\mathbb{N} \to \text{List } \mathbb{N}$ nats = from zerofrom n = n :: from (succ n)

However, in our setting, strong normalisation must be preserved to ensure logical consistency. The definition of from cannot be accepted as it performs a recursive call on a term which is not structurally smaller than the argument.

We would still like to represent such objects, but care must be taken to construct them in a coherent way. The class of *coinductive types* [13,19] is introduced to this end. For example, the **Stream** coinductive family collects the types of infinite sequences.

> coinductive Stream : Type  $\rightarrow$  Type  $\triangleright: \Pi_{(A:\mathsf{Type})} A \rightarrow \mathsf{Stream} A \rightarrow \mathsf{Stream} A$

The definition is very similar to that of the List inductive family. We omit an empty constructor here, so these streams are necessarily infinite, though an alternative definition of lists which are potentially infinite is also possible, which more closely aligns with lazy lists as they exist in Haskell.

Crucially, declaring the type coinductive lifts the requirement that terms of this type must be wellfounded. If we had instead declared the type inductive, it could be shown that Stream  $A \simeq 0$ , as we could never construct a base value for the constructor's second argument. In particular, we don't require functions constructing streams to be structurally recursive, so the above definitions can be ported to result in Stream  $\mathbb{N}$ .

 $\begin{array}{ll} \mathsf{nats}:\mathsf{Stream}\ \mathbb{N} & \quad \mathsf{from}:\mathbb{N}\to\mathsf{Stream}\ \mathbb{N}\\ \mathsf{nats}=\mathsf{from}\ \mathsf{zero} & \quad \mathsf{from}\ n=n \mathrel{\triangleright}\mathsf{from}\ (\mathsf{succ}\ n) \end{array}$ 

However, not all recursive definitions are allowed. We must ensure canonicity, so that a term of type **Stream** A is actually constructed by  $\triangleright$  applied to the appropriate arguments. This will justify being able define functions by pattern matching on the type, such as **head** and **tail** projection functions.

$$\begin{array}{ll} \mathsf{head}: \Pi_{(A:\mathsf{Type})} \ \mathsf{Stream} \ A \to A & \mathsf{tail}: \Pi_{(A:\mathsf{Type})} \ \mathsf{Stream} \ A \to \mathsf{Stream} \ A \\ \mathsf{head} \ (h \triangleright t) = h & \mathsf{tail} \ (h \triangleright t) = t \end{array}$$

Strong normalisation is preserved if we only reduce recursive calls lazily [25, 32]. That is, only when the result of reducing the term is needed in order to reduce a pattern-matching function.

In general, recursive calls must be *guarded* by a constructor of the coinductive type [28, 47], that is, occurring within a subexpression of a constructor. This guarantees the function must make some progress at each recursive step. In the definition of from, we can see that guardedness is satisfied, the recursive call appearing as the second argument of the  $\triangleright$  constructor.

An example that is ruled out by the guardedness condition is the function which filters a stream according to a decidable predicate, keeping only those elements which satisfy the predicate. For lists, this has a valid definition by structural recursion.

filter : 
$$\Pi_{(A:\mathsf{Type})} (A \to 2) \to \mathsf{List} A \to \mathsf{List} A$$
  
filter  $f [] = []$   
filter  $f (h :: t) = \mathsf{elim}_+ (\mathsf{elim}_1 (h :: \mathsf{filter} f t)) (\mathsf{elim}_1 (\mathsf{filter} f t)) (f h)$ 

This definition has two recursive calls on the tail of the list. The first is when f h is true, so h will be kept in the resulting list, h :: filter f t, while f h is false in the second, dropping h from the result, filter f t.

When attempting to define a similar function for type Stream A, the second recursive call would be problematic as it makes no progress in constructing a new  $h \triangleright t$  node. In general, there is no guarantee that the predicate will ever return true — given a stream s: Stream A, the term head (filter ( $\lambda(a : A)$ . false) s) would loop forever, never producing a head for the new stream.

Other familiar list functions may have corresponding definitions for streams as long as they are productive. For example, a function may be mapped pointwise over an stream to yield a new stream, and two streams may be zipped to form a stream of pairs.

$$\begin{split} \mathsf{map}: \Pi_{(A,B:\mathsf{Type})} \ (A \to B) \to \mathsf{Stream} \ A \to \mathsf{Stream} \ B \\ \mathsf{map} \ f \ (h \triangleright t) = f \ h \triangleright \mathsf{map} \ f \ t \end{split}$$

 $\begin{aligned} \mathsf{zip} : \Pi_{(A,B:\mathsf{Type})} \text{ Stream } A &\to \mathsf{Stream } B \to \mathsf{Stream } (A \times B) \\ \mathsf{zip} \ (h_1 \triangleright t_1) \ (h_2 \triangleright t_2) &= \langle h_1, \ h_2 \rangle \triangleright \mathsf{zip } \ t_1 \ t_2 \end{aligned}$ 

Coinductive types may also represent propositions whose proofs are infinite structures. For example, the Always P predicate on streams asserts that when predicate P holds at the head of a stream and, recursively, Always P holds on the tail, then it holds at every element.

$$\begin{array}{l} \mbox{coinductive Always}: \Pi_{(A: {\sf Type})} \ (A \to {\sf Type}) \to {\sf Stream} \ A \to {\sf Type} \\ \mbox{always}: \Pi_{(P:A \to {\sf Type})} \ \Pi_{(s: {\sf Stream} \ A)} \\ P \ ({\sf head} \ s) \to {\sf Always} \ P \ ({\sf tail} \ s) \to {\sf Always} \ P \ s \end{array}$$

A proof of Always  $P \ s$  can itself be seen as an infinite stream of proofs, albeit each element having a different type, the first being P (head s), the second P (head (tail s)), then P (head (tail s))), and so on.

Two streams are extensionally equal when they are pointwise equal. Even if

their definitions may radically differ, they are considered equivalent if they produce the same sequence of values. This form of equivalence is called *bisimulation*, and can be defined in terms of the Always type.

Bisim : 
$$\Pi_{(A:\mathsf{Type})}$$
 Stream  $A \to \mathsf{Stream} \ A \to \mathsf{Type}$   
Bisim  $r \ s = \mathsf{Always} \ (\lambda(p : A \times A). \ \pi_1 \ p \equiv \pi_2 \ p) \ (\mathsf{zip} \ r \ s)$ 

Or, equivalently, as a dedicated coinductive type, a binary type family indexed by streams.

coinductive Bisim :  $\Pi_{(A:\mathsf{Type})}$  Stream  $A \to \mathsf{Stream} \ A \to \mathsf{Type}$ bisim :  $\Pi_{(a,b:\mathsf{Stream} \ A)}$  head  $a \equiv \mathsf{head} \ b \to \mathsf{Bisim} \ (\mathsf{tail} \ a) \ (\mathsf{tail} \ b) \to \mathsf{Bisim} \ a \ b$ 

For example, we can show that mapping succ onto stream from n is the same as applying tail. These are two streams that can be proved bisimilar by coinduction.

**Theorem 2.7.1.**  $\Pi_{(n:\mathbb{N})}$  Bisim (map succ (from n)) (tail (from n))

*Proof.* First, use the bisim constructor to form a term of the coinductive Bisim type. We now need to provide two proofs, one showing the equality of the streams' heads, the other showing that their tails are bisimilar.

For the first proof, observe that the heads of both streams reduce to succ n, so they may be proven equal by  $refl_{\equiv}$  (succ n).

head (map succ (from n)) = succ (head (from n)) (definition of map) = succ n (definition of from)

```
 \begin{array}{ll} \mbox{head (tail (from <math>n)) = head (from (succ n)) (definition of from) \\ = \mbox{succ } n & (definition of from) \end{array}
```

For the bisimilarity of the tails, the relaxed limits of recursion allow us to assume the statement we are proving as a coinductive hypothesis. Since we are guarded by an application of constructor **bisim**, we may now use the hypothesis without further restriction.

Instantiating the coinductive hypothesis at index succ n gives the statement Bisim (map succ (from (succ n))) (tail (from (succ n))). It can be seen that that this is judgementally equal to the statement we want to prove.

 $map \ succ \ (from \ (succ \ n)) = map \ succ \ (tail \ (from \ n)) \qquad (definition \ of \ from)$  $= tail \ (map \ succ \ (from \ n)) \qquad (definition \ of \ map)$ 

tail (from (succ n)) = tail (tail (from n)) (definition of from)

So we have Bisim (tail (map succ (from n))) (tail (tail (from n))) by coinductive hypothesis, which is precisely the statement that the tails of the two streams are bisimilar.

Formally, the proof can be written as a function defined by the following recursive equation, where the recursive call is guarded by the **bisim** constructor.

 $f: \Pi_{(n:\mathbb{N})} \text{ Bisim (map succ (from n)) (tail (from n))}$  $f n = \text{bisim (refl}_{\equiv} (\text{succ } n)) (f (\text{succ } n))$ 

Applying the function to 0 gives a proof that map succ nats and tail nats are bisimilar. Unravelling the recursive calls in the proof term, we see that it actually an infinite sequence of  $refl_{\equiv}$  proofs, a new one generated at each recursive step.

$$\begin{split} f & 0 = \text{bisim } (\text{refl}_{\pm} \ 1) \ (f \ 1) \\ & = \text{bisim } (\text{refl}_{\pm} \ 1) \ (\text{bisim } (\text{refl}_{\pm} \ 2) \ (f \ 2)) \\ & = \text{bisim } (\text{refl}_{\pm} \ 1) \ (\text{bisim } (\text{refl}_{\pm} \ 2) \ (\text{bisim } (\text{refl}_{\pm} \ 3) \ (f \ 3))) \\ & = \text{bisim } (\text{refl}_{\pm} \ 1) \ (\text{bisim } (\text{refl}_{\pm} \ 2) \ (\text{bisim } (\text{refl}_{\pm} \ 3) \ (\dots))) \end{split}$$

Bisimilarity is an equivalence relation. The properties of reflexivity, symmetry, and transitivity may each be shown by coinductive proof. We demonstrate these by giving a recursive function for each.

Theorem 2.7.2. Reflexivity, Bisim a a.

Proof.

$$\operatorname{refl}_{\mathsf{Bisim}} : \Pi_{(A:\mathsf{Type})} \Pi_{(a:\mathsf{Stream}\ A)} \operatorname{Bisim} a \ a$$
$$\operatorname{refl}_{\mathsf{Bisim}} (h \triangleright t) = \operatorname{bisim} (\operatorname{refl}_{\equiv} h) (\operatorname{refl}_{\mathsf{Bisim}} t)$$

We define the function by pattern matching on the stream, then use the bisim constructor to form the resulting proof. We must therefore prove  $h \equiv h$  and Bisim t t.

Using the reflexivity constructor of propositional equality, the head is equal to itself,  $\operatorname{refl}_{\equiv} h$ . To prove the tail bisimilar to itself, since we are guarded by an application of bisim, we may recursively call the function at the tail, proving the statement by coinductive hypothesis.

**Theorem 2.7.3.** Symmetry, Bisim  $a \ b \rightarrow$  Bisim  $b \ a$ .

Proof.

$$\begin{array}{l} \mathsf{sym}_{\mathsf{Bisim}}: \Pi_{(A:\mathsf{Type})} \ \Pi_{(a,b:\mathsf{Stream}\ A)} \ \mathsf{Bisim}\ a\ b \to \mathsf{Bisim}\ b\ a\\ \mathsf{sym}_{\mathsf{Bisim}}\ (\mathsf{bisim}\ h\ t) = \mathsf{bisim}\ (\mathsf{sym}_{\equiv}\ h)\ (\mathsf{sym}_{\mathsf{Bisim}}\ t) \end{array}$$

Here, we pattern match on the proof of Bisim  $a \ b$ , which is formed by proofs  $h : \text{head } a \equiv \text{head } b$  and t : Bisim (tail a) (tail b). After using the bisim constructor, we must prove the symmetric forms of these assumptions.

The first, head  $b \equiv \text{head } a$ , is clear from the symmetry of equality,  $\text{sym}_{\equiv} h$ . The second, Bisim (tail b) (tail a), can be shown by coinductive hypothesis, recursively calling the function on t.

**Theorem 2.7.4.** Transitivity, Bisim  $a \ b \rightarrow Bisim \ b \ c \rightarrow Bisim \ a \ c$ .

Proof.

 $\begin{array}{l} \operatorname{trans}_{\mathsf{Bisim}}: \Pi_{(A:\mathsf{Type})} \ \Pi_{(a,b,c:\mathsf{Stream}\ A)} \ \mathsf{Bisim}\ a\ b \to \mathsf{Bisim}\ b\ c \to \mathsf{Bisim}\ a\ c \\ \operatorname{trans}_{\mathsf{Bisim}}\ (\mathsf{bisim}\ h_1\ t_1)\ (\mathsf{bisim}\ h_2\ t_2) = \mathsf{bisim}\ (\mathsf{trans}_{\equiv}\ h_1\ h_2)\ (\mathsf{trans}_{\mathsf{Bisim}}\ t_1\ t_2) \\ \end{array}$ 

The structure of this proof is the same as that of symmetry, pattern matching on the proofs of bisimilarity. We use the transitive property of equality for the head, and the coinductive hypothesis for the tail.  $\Box$ 

As a final note for this introduction, we have focused on infinite linear sequences as these are most relevant for Chapter IV, but coinductive types for other kinds of non-wellfounded tree may also be defined. A type of infinitely-deep binary trees is produced with a constructor that takes two recursive arguments instead of one.

> coinductive Tree : Type  $\rightarrow$  Type node :  $\Pi_{(A:\mathsf{Type})} A \rightarrow \mathsf{Tree} A \rightarrow \mathsf{Tree} A \rightarrow \mathsf{Tree} A$

The requirement for guarded recursion remains the same. For example, the binary tree with whose nodes state their index in a breadth-first traversal can be defined by the following equation whose recursive calls are guarded by the node constructor.

$$\begin{split} & \mathsf{bft}: \mathbb{N} \to \mathsf{Tree} \ \mathbb{N} \\ & \mathsf{bft} \ n = \mathsf{node} \ n \ (\mathsf{bft} \ (\mathsf{add} \ 1 \ (\mathsf{add} \ n \ n))) \ (\mathsf{bft} \ (\mathsf{add} \ 2 \ (\mathsf{add} \ n \ n))) \end{split}$$

Applying the **bft** function to argument 0 gives the actual tree.

```
bft \ 0 = node \ 0 \ (node \ 1 \ (node \ 3 \ (...) \ (...)) \ (node \ 4 \ (...) \ (...)))(node \ 2 \ (node \ 5 \ (...) \ (...)) \ (node \ 6 \ (...) \ (...)))
```



If val, left, and right are the projection functions for the node constructor, tree bisimulation may then be defined as a coinductive family indexed by trees.

coinductive  $\operatorname{Bisim}_{\operatorname{Tree}} : \Pi_{(A:\operatorname{Type})}$  Tree  $A \to \operatorname{Tree} A \to \operatorname{Type}$  $\operatorname{bisim}_{\operatorname{Tree}} : \Pi_{(a,b:\operatorname{Tree} A)}$  val  $a \equiv \operatorname{val} b \to \operatorname{Bisim}_{\operatorname{Tree}}$  (left a) (left b)  $\to$  $\operatorname{Bisim}_{\operatorname{Tree}}$  (right a) (right b)  $\to$   $\operatorname{Bisim}_{\operatorname{Tree}} a b$ 

Its status as an equivalence relation can be proved analogously to bisimilarity of streams, by coinduction. For each property, the defining recursive function will have two recursive calls, instantiated separately for the left and right subtrees.

# Chapter III Epistemic Modal Logic

In this chapter, we will introduce the axiomatic system of epistemic modal logic and its standard relational semantics. Initially, this will be presented in isolation from the system of type theory defined Chapter II, but we will later go on to define a shallow embedding of the logic into our type theory.

The embedding constitutes a novel semantics for epistemic logic based on *knowledge operators*, which will be shown to correspond to the relational semantics. An Agda formalisation of the embedding and the semantics correspondence results are included in Appendix A.2.

# 3.1 Background

Modal logics are a rich field of study throughout mathematics and philosophy, as well as computer science. They were created in their modern form to study the linguistic modality "necessarily", which qualifies the truth of a statement [20, 27, 40]. For example, the statements "it is raining" and "it is necessarily raining" have different denotations,

Formally, the syntax of a propositional logic is extended with a new unary operator  $\Box$  on formulae, with  $\Box \phi$  interpreted as the statement " $\phi$  is necessary". There are a number of axiomatisations for the new operator, the weakest of which is a logic called *System K*, which introduces two new rules. The rules use the judgement  $\phi_1, \ldots, \phi_n \vdash \psi$  which asserts that the propositional formula  $\psi$  is derivable under hypotheses  $\phi_1$  through  $\phi_n$ . When  $\psi$  can be derived axiomatically, we may write  $\vdash \psi$ . Note that while the notation is similar, these rules are defining a distinct logical system from the type theory defined in Chapter II.

NECESSITATION 
$$\frac{\vdash \phi}{\vdash \Box \phi}$$
 AXIOM K  $\overline{\vdash \Box (\phi \Rightarrow \psi) \Rightarrow \Box \phi \Rightarrow \Box \psi}$ 

The necessitation rule states that when  $\phi$  is a tautology of the underlying propositional logic, then so is  $\Box \phi$ . That is, tautologies of the logic are necessary.

Axiom K asserts necessity distributes over implication. It is a version of the modus ponens rule that can be applied to necessary formulae — if both a formula and an implication with that formula as the premise are necessary, then the conclusion must also be necessary.

To strengthen the logic, several more axioms may be added depending on the intended use. We will skip straight to the modal logic *System 5 (S5)* including the necessitation rule, Axiom K, and three additional axioms.

AXIOM T 
$$\overline{+ \Box \phi \Rightarrow \phi}$$
 AXIOM 4  $\overline{+ \Box \phi \Rightarrow \Box \Box \phi}$   
AXIOM 5  $\overline{+ \neg \Box \phi \Rightarrow \Box \neg \Box \phi}$ 

These axioms dictate what can be derived from the assumption that  $\phi$  is necessary or unnecessary. From Axioms T and 4, if a formula is necessary, then it must actually be true, and its necessity is necessary.

In some presentations of modal logic, Axiom 5 is written  $\Diamond \phi \Rightarrow \Box \Diamond \phi$ . This uses a new modal operator  $\Diamond$ , with  $\Diamond \phi$  meaning " $\phi$  is possible", with the axiom reading that if a formula is possible, it is necessarily possible. Possibility may be defined in terms of necessity,  $\Diamond \phi = \neg \Box \neg \phi$ , that is, a formula is possible when it's not necessarily false. The above version of the axiom, with a single modal operator, is classically equivalent and more suited for use in epistemic modal logic.

If tautological formulae are necessary, and necessary formulae are true, it might not be clear what we gain by adding the  $\Box$  modality — what is the difference between the propositions  $\Box \phi$  and  $\phi$ ? The answer is that modal logics are used to model scenarios containing unknowns, where we may be in any one of a number of *possible worlds*, where some formulae might be true, but cannot be derived tautologically.

For example, our earlier statement "it is raining" is true on some days, but false on others. In a mathematical game with hidden information, such as card games or a prisoner dilemma, the state of the play may not be known by a particular player — the truth of the statement "my opponent holds the  $Q \heartsuit$ card" is contingent on the state. We add primitive formulae, called *events*, to the logic depending on the scenario, whose truth value is a function of the precise state of the world.

The traditional reading of  $\Box \phi$  as " $\phi$  is necessary" can be modified, with different axioms selected, according to the scenario. In philosophy, moral obligation may be studied with a *deontic modal logic*, where  $\Box \phi$  is read " $\phi$  is obligatory", with  $\Diamond \phi$  therefore meaning " $\phi$  is permitted", and which lacks Axiom T — even if an action is obligatory, it may not actually occur. In a *temporal modal logic*, the modal operator refers to tense, one possible reading of  $\Box \phi$  being " $\phi$  is true at all times", and  $\Diamond \phi$  meaning " $\phi$  is true some of the time". In this chapter, we are concerned with *epistemic modal logic* [24, 33], where the modal operator is written K  $\phi$  with the interpretation " $\phi$  is known".

# **3.2** Epistemic Logic and Relational Semantics

In this section, we will see how the S5 modal logic forms an epistemic logic, appropriate for the modelling of knowledge. We will also discuss the traditional semantics of modal logic, based on binary relations, and how these relations are interpreted in epistemic logic.

We have a modal operator which can be applied to logical formulae, K  $\phi$  interpreted as " $\phi$  is known", but known by whom? It could be a player in a mathematical game, a computer on a network, or a Byzantine general planning a coordinated attack on an enemy city [38]. We will not presume any particular scenario, leaving the subject as an abstract *agent*. A natural extension is postulating many such agents, each agent *a* having their own modal operator  $K_a$ , but for now we will only consider a logic with one K operator.

The S5 modal logic introduced in Section 3.1 can be utilised to give an account of an agent's knowledge. Let us consider each rule, giving a brief justification for how it is interpreted in epistemic logic.

• KNOWLEDGE GENERALISATION, if  $\vdash \phi$ , then  $\vdash \mathsf{K} \phi$ 

This is the necessitation rule of System K, called *knowledge generali*sation in the context of epistemic logic. It states that if  $\phi$  is derivable, then so is K  $\phi$ . In other words, the agent is capable of using logic to derive theorems.

This may seem an overly-strong assumption, tantamount to the agent being a perfect logician. However, consider that in order for us to derive  $\mathsf{K} \ \phi$  by knowledge generalisation, we must actually have a proof of  $\phi$ . If formula  $\psi$  is the statement of an unsolved problem in mathematics, such as the Collatz conjecture, at present we can show neither  $\mathsf{K} \ \psi$  nor  $\mathsf{K} \neg \psi$ .

• AXIOM K, K  $(\phi \Rightarrow \psi) \Rightarrow$  K  $\phi \Rightarrow$  K  $\psi$ 

If both an implication and its premise are known, then so is its conclusion. That is, the agent is capable of using modus ponens with their own knowledge.

In cases where both  $\phi$  and  $\phi \Rightarrow \psi$  are theorems, then so is  $\psi$ , and knowledge generalisation is sufficient to conclude K  $\psi$ . However, this rule

allows for the agent to apply logic to known events which are not derivable tautologically.

• AXIOM T, K  $\phi \Rightarrow \phi$ 

If a formula is known, it must actually be true in the present state of the world. This is another strong assumption, but this is what distinguishes knowledge from mere belief or opinion.

Belief can itself be expressed in *doxastic modal logic*, which might adopt a weaker axiom of belief non-contradiction, that an agent cannot simultaneously believe a statement and its negation,  $\neg(\mathsf{B} \phi \land \mathsf{B} \neg \phi)$ .

• AXIOM 4,  $\mathsf{K} \phi \Rightarrow \mathsf{K} (\mathsf{K} \phi)$ 

If a formula is known, then this knowledge is itself known. This is a principle of self-reflection, that our agent is aware of their own knowledge.

• AXIOM 5,  $\neg \mathsf{K} \phi \Rightarrow \mathsf{K} (\neg \mathsf{K} \phi)$ 

If a formula isn't known, then this lack of knowledge is itself known. This is a negative version of Axiom 4, when the agent doesn't know something, they at least know they do not know it.

When these two axioms are taken together, the agent is perfectly introspective, aware of both what they know and what they don't. This could be called a Socratic principle — "I know that I know nothing".

The relational semantics of modal logic, and therefore epistemic logic, was originally devised by Saul Kripke, who proved it sound and complete with respect to the axiomatic system [37]. The semantics operates on a set of possible worlds, relating them according to the knowledge of the agent.

Formally, a Kripke frame is a pair  $\langle S, R \rangle$  where S is a set of world states and R is a binary relation on states. A Kripke frame with an arbitrary relation R is sufficient to model the modal logic System K, but additional properties can be imposed on the relation which correspond to additional axioms. The following the properties are required of R to model the modal logic S5.

- KNOWLEDGE GENERALISATION, if  $\vdash \phi$ , then  $\vdash \mathsf{K} \phi$ Modelled by any relation R.
- AXIOM K,  $\mathsf{K} \ (\phi \Rightarrow \psi) \Rightarrow \mathsf{K} \ \phi \Rightarrow \mathsf{K} \ \psi$ Modelled by any relation R.
- AXIOM T,  $\mathsf{K} \phi \Rightarrow \phi$ Modelled when R is reflexive,  $w \mathrel{R} w$ .
- AXIOM 4,  $\mathsf{K} \phi \Rightarrow \mathsf{K} (\mathsf{K} \phi)$ Modelled when *R* is transitive,  $w \ R \ v \Rightarrow v \ R \ u \Rightarrow w \ R \ u$ .

• AXIOM 5,  $\neg \mathsf{K} \phi \Rightarrow \mathsf{K} (\neg \mathsf{K} \phi)$ Modelled when *R* is Euclidean,  $w \mathrel{R} v \Rightarrow w \mathrel{R} u \Rightarrow u \mathrel{R} v$ .

A reflexive, Euclidean relation must also be symmetric, and a symmetric, transitive relation is Euclidean, so S5 is modelled precisely when R is an equivalence relation. We omit the full details of the correspondence here as a version of it, embedded in type theory, is formally proven in Section 3.5.

Intuitively, the relation, called an *accessibility relation* or *knowledge relation* in the context of epistemic logic, relates the states of the world that agent's knowledge cannot distinguish. That is, two states w and v are related, w R v, when they differ only by the truth of those formulae  $\phi$  for which  $\neg \mathsf{K} \phi$ .

If  $\forall_{(w\in S)} \exists_{(v\in S)} w \ R \ v \Rightarrow w = v$ , that is, each state is disconnected from all others, beside itself reflexively, the agent modelled by R is omniscient, able to precisely determine the state of the world. If the agent were then to forget an event  $\phi$ , the relation would expand to form equivalence classes relating those states which only differ by the truth of  $\phi$ . In the extreme case, when  $\forall_{(w\in S)}\forall_{(v\in S)} w \ R \ v$ , when R is strongly connected, the agent is solipsistic, completely ignorant of the world around them, only knowing tautologies through knowledge generalisation and the state of their own knowledge though Axioms 4 and 5.

# 3.3 Type-Theoretic Embedding

We now return to a type-theoretic setting to define an embedding of epistemic modal logic into the logic of type theory. Broadly, there are two methods to embed one logic into another [51], serving different purposes depending on the desired application.

In a *deep embedding*, the syntax of the object logic is defined directly in the metalogic. For example, a deep embedding of epistemic modal logic in type theory could utilise an inductive type to define its syntax trees.

 $\begin{array}{l} \mathsf{inductive \ Syntax: Type} \\ \bot : \mathsf{Syntax} \\ \neg : \mathsf{Syntax} \to \mathsf{Syntax} \\ \lor : \mathsf{Syntax} \to \mathsf{Syntax} \to \mathsf{Syntax} \\ \mathsf{K} : \mathsf{Syntax} \to \mathsf{Syntax} \end{array}$ 

To keep the eliminator for the inductive type as simple as possible, only the essential connectives are defined as its constructors. The others may be defined in terms of these using standard classical equivalences.

op : Syntax	$\wedge:Syntax\toSyntax\toSyntax$	
$\top = \neg \bot$	$\phi \wedge \psi =$	$= \neg (\neg \phi \lor \neg \psi)$
$\Rightarrow:Syntax\toSyntax\to$	Syntax	$\Leftrightarrow:Syntax\toSyntax\toSyntax$
$\phi \Rightarrow \psi = \neg \phi \lor \psi$		$\phi \Leftrightarrow \psi = (\phi \Rightarrow \psi) \land (\psi \Rightarrow \phi)$

An axiomatic system is specified over the syntax by defining a relation to establish the valid derivations of the logic. Here, an inductive type family indexed by both Syntax and List Syntax is used, the list representing prior assumptions.

inductive Derives : List Syntax  $\rightarrow$  Syntax  $\rightarrow$  Type

K-gen	: $\Pi_{(\phi:Syntax)}$	$Derives \ [ \ ] \ \phi \to Derives \ [ \ ] \ (K \ \phi)$
$axiom_T$	: $\Pi_{(\phi:Syntax)}$	Derives [] (K $\phi \Rightarrow \phi$ )
$axiom_K$	: $\Pi_{(\phi,\psi:Syntax)}$	$Derives\;[\;]\;(K\;(\phi\Rightarrow\psi)\RightarrowK\;\phi\RightarrowK\;\psi)$
$axiom_4$	: $\Pi_{(\phi:Syntax)}$	$Derives [ ] (K \phi \Rightarrow K (K \phi))$
$axiom_5$	: $\Pi_{(\phi:Syntax)}$	$Derives [] (\neg K \phi \Rightarrow K (\neg K \phi))$

A closed term of type Derives  $\Delta \phi$  is a proof tree for the derivation of  $\Delta \vdash \phi$ in epistemic logic. Only the additional rules of epistemic logic are shown above, but all of the rules for the base propositional logic will also need be included in the relation.

This style of embedding is appropriate for studying the metatheoretic properties of the logic [39]. An interpretation function may be defined, mapping from Syntax to a semantic model, with soundness and completeness results able to be stated and proved in the metalogic of type theory by structural induction. However, proving theorems of the object logic using a deep embedding can be cumbersome — the rules for both the object logic and the metalogic must be observed.

Alternatively, a *shallow embedding* may be used if the objective is theorem proving in the object logic. The metalogic is used as a semantic model of the object logic — rather than confining syntax trees to a single type, the connectives of the metalogic are used to encode them. This allows for the full feature set of the metalogic to be brought to bear in embedded proofs [12].

The rest of this chapter is concerned with a shallow embedding of epistemic logic in type theory. We will see how type-theoretic features such as type families and, in Chapter IV, coinduction may be used to define the rules of epistemic logic at a higher level than is possible with a deep embedding.

The embedding is based on the concept of possible worlds, so we fix a type

to represent these worlds, State : Type. A formula in the syntax of epistemic logic may be interpreted as true in some worlds but false in others, so embedded formulae are events, predicates on states, encoded as type families indexed by State.

**Definition 3.3.1.** An event is a type family indexed by the type of possible worlds.

 $\begin{array}{l} \mathsf{Event}:\mathsf{Type}\\ \mathsf{Event}=\mathsf{State}\to\mathsf{Type} \end{array}$ 

As a reminder, to ensure consistency, since both State and Type are used in its definition, there are constraints on which type universe Event itself inhabits. Its level must be at least as high as the universe of State and strictly higher than the universe that is the codomain of the type family, with the actual level being the least value which satisfies both constraints. With this noted, we shall continue being implicit about levels in future definitions — the Agda formalisation in Appendix A details them in full.

The syntax of a propositional logic is built atop the constructs of type theory. The constants are constant type families, while the connectives are functions combining type families using the type formers introduced in Chapter II.

Definition 3.3.2. Event constants and connectives.

$\perp$ : Event	op: Event		$\neg:Event\toEvent$	
$\perp = \lambda(w : State). \ \mathbb{O}$	$ op = \lambda(w : State). \ 1$		$\neg \phi = \lambda(w: State). \ \phi \ w \to \mathbb{O}$	
$\vee:Event\toEvent-$	> Event	$\wedge:E$	$vent \to Event \to Event$	
$\phi \lor \psi = \lambda(w:State)$	. $\phi w + \psi w$	$\phi \wedge \chi$	$\psi = \lambda(w : State). \ \phi \ w  imes \psi \ w$	
$\Rightarrow:Event\toEvent\toI$	Event	⇔: E	Event $ ightarrow$ Event $ ightarrow$ Event	
$\phi \Rightarrow \psi = \lambda(w : State). \ \phi \ w \rightarrow \psi \ w$		$\phi \Leftrightarrow \psi = \lambda(w : State). \ \phi \ w \leftrightarrow \psi \ w$		

An event may be seen set-theoretically as the subset of states where the event occurs. With this view, applying an event to a state,  $\phi w$ , is obtaining the truth assignment of  $w \in \phi$ . The constants  $\perp$  and  $\top$  are the empty and universal sets,  $\neg$  is the complement, with  $\lor$ ,  $\land$ , and  $\Rightarrow$  being the union, intersection, and exponent, respectively.

We can also define predicates and relations over events. Rather than representing connectives in the syntax of the object logic, these correspond to metatheoretic statements — a shallow embedding allows these to coexist in the same language. Definition 3.3.3. Event predicates and relations.

$$\begin{split} &\forall: \mathsf{Event} \to \mathsf{Type} \\ &\forall \phi = \Pi_{(w:\mathsf{State})} \ \phi \ w \\ &\subseteq: \mathsf{Event} \to \mathsf{Event} \to \mathsf{Type} \\ &\phi \subseteq \psi = \forall (\phi \Rightarrow \psi) \\ &\phi \approx \psi = \forall (\phi \Leftrightarrow \psi) \end{split}$$

The statement  $\forall \phi$  states that  $\phi$  is a tautology in the object logic, true in every possible world. Continuing the set theory analogy, it means that  $\phi$  is the universal set, with  $w \in \phi$  for any state w. The tautological implication,  $\phi \subseteq \psi$ , is the subset relation, whenever  $w \in \phi$ , then  $w \in \psi$ . Two events are propositionally equivalent,  $\phi \approx \psi$ , when this holds bidirectionally, being equal sets.

One advantage of the shallow embedding style is immediately apparent. As we have defined the connectives of the embedded propositional logic using those of type theory, their properties follow straightforwardly with little extra effort. For example, unwrapping the definitions in  $\phi \wedge \psi \subseteq \psi \wedge \phi$ , it can be seen to be nothing more than the type-theoretic theorem that products are commutative, specialised to type families  $\phi$  and  $\psi$ .

$$\Pi_{(w:\mathsf{State})} \phi \ w \times \psi \ w \to \psi \ w \times \phi \ w$$

This embedding is enough to begin modelling some basic scenarios in epistemic logic. For example, imagine that a coin has been tossed and a six-sided die rolled. Our chosen type of states should at least distinguish these primitive events.

		$D_1$	"the die rolled a 1"
$C_H$	"the coin landed heads-side up"	1	•
$C_T$	"the coin landed tails-side up"		:
- 1		$D_6$	"the die rolled a 6"

Then, for example,  $D_2 \vee D_4 \vee D_6$  is the event that occurred in those states where the die rolled an even number. We would also expect certain statements like  $\forall \neg (C_H \wedge C_T)$  to hold, so there cannot be inconsistent states where the coin or die give multiple outcomes. One possible choice of states for this scenario is State = Fin 2 × Fin 6, using finite types to represent the outcomes for the coin and die independently.

Notably absent, however, are any knowledge modalities. Being unary operators in the syntax of the embedded epistemic logic, it is clear that they should have type  $\mathsf{Event} \to \mathsf{Event}$ . Their semantics will be defined formally in Section 3.4, but for now assume we do have a function  $\mathsf{K}_a$  representing the knowledge of some agent a. In our example, we now have additional events to model,  $\mathsf{K}_a D_1$ ,  $\mathsf{K}_a D_2$ , and so on, but also more complex events like  $\neg(\mathsf{K}_a C_H \lor \mathsf{K}_a C_T)$ , "agent a doesn't know on which side the coin landed". The state type Fin  $2 \times$  Fin 6 is no longer rich enough to encode this detail.

If we then include another agent b with associated knowledge operator  $\mathsf{K}_b$ , we also have events about one agents knowledge of the other's knowledge to contend with. For example, the event  $\mathsf{K}_a$  ( $\mathsf{K}_b$   $D_3$ )  $\land \neg \mathsf{K}_b$  ( $\mathsf{K}_a$  ( $\mathsf{K}_b$   $D_3$ )) states "agent a knows that agent b knows the die rolled a 3, but b doesn't know that aknows this". Due to this added complexity, it is simpler to only talk abstractly about states, postulating properties such as the coin consistency principle stated above, rather than giving **State** a concrete definition which must distinguish the knowledge of agents.

# 3.4 Knowledge Operator Semantics

We now come to defining what it means to be a knowledge operator in our embedding. Instead of defining a function  $K : Event \rightarrow Event$  directly as we have with the other syntactic constructs of epistemic logic, we shall define them semantically, stating a set of properties a function must satisfy in order to be appropriate to represent the knowledge of an agent. The semantics will be based on the rules of the S5 system of modal logic introduced in Sections 3.1 and 3.2, defined in the syntax of our embedded epistemic logic.

$\forall\!\!\!\!\forall \phi \to \forall\!\!\!\!\forall K \phi$	(Knowledge Generalisation)
$K \ (\phi \Rightarrow \psi) \subseteq K \ \phi \Rightarrow K \ \psi$	(Axiom K)
$K \ \phi \subseteq \phi$	(Axiom T)
$K \ \phi \subseteq K \ (K \ \phi)$	(Axiom 4)
$\neg K \ \phi \subseteq K \ (\neg K \ \phi)$	(Axiom 5)

We are not limited to the finitary syntax of the object logic, however. We are able to make stronger semantic statements which may involve entailments to and from potentially infinite sets of formulae.

Recall that a type family indexed by type T is a function  $T \to \mathsf{Type}$  which may be thought of as a set of types. In the same way, we define event families as functions  $T \to \mathsf{Event}$ , which are sets of formulae in epistemic logic. In full, an event family is a function  $T \to \mathsf{State} \to \mathsf{Type}$ , a binary type family indexed by Tand  $\mathsf{State}$ .

**Definition 3.4.1.** An event family  $\Phi$  semantically entails an event family  $\Psi$  when, in every state, if  $\Phi$  holds universally, at all of its indices, then  $\Psi$  must also

hold universally.

$$\models : (A \to \mathsf{Event}) \to (B \to \mathsf{Event}) \to \mathsf{Type}$$
  
 
$$\Phi \models \Psi = \Pi_{(w:\mathsf{State})} \ (\Pi_{(a:A)} \ \Phi \ a \ w) \to (\Pi_{(b:B)} \ \Psi \ b \ w)$$

By composing an event family  $\Phi$  with a function  $\mathsf{K} : \mathsf{Event} \to \mathsf{Event}$ , we form a new event family, indexed by the same type, where  $\mathsf{K}$  is mapped onto the family. When  $\Phi \vDash \Psi$  implies  $\mathsf{K} \circ \Phi \vDash \mathsf{K} \circ \Psi$ , we say that  $\mathsf{K}$  preserves semantic entailment. We require that knowledge operators preserve semantic entailment in place of the knowledge generalisation and the Axiom K property.

In what follows, we are most often concerned about entailments from an event family into a single event, instantiating  $\Psi$  with a constant event family,  $\lambda b$ .  $\psi$ . We shall write this without the extraneous  $\lambda$ -term,  $\Phi \models \psi$ .

**Definition 3.4.2.** A knowledge operator is any function  $K : Event \rightarrow Event$  satisfying preservation of semantic entailment in addition to the S5 properties Axiom T, Axiom 4, and Axiom 5. That is, the type of knowledge operators is defined as the following dependent tuple type.

$$\begin{split} \mathsf{KOp}:\mathsf{Type} \\ \mathsf{KOp} &= \Sigma_{(\mathsf{K}:\mathsf{Event}\to\mathsf{Event})} & \Pi_{(\Phi:\mathsf{Event}\to\mathsf{Type})} \Pi_{(\psi:\mathsf{Event})} & \Phi \vDash \psi \to \mathsf{K} \circ \Phi \vDash \mathsf{K} \ \psi \\ & \times & \Pi_{(\phi:\mathsf{Event})} & \mathsf{K} \ \phi \subseteq \phi \\ & \times & \Pi_{(\phi:\mathsf{Event})} & \mathsf{K} \ \phi \subseteq \mathsf{K} \ (\mathsf{K} \ \phi) \\ & \times & \Pi_{(\phi:\mathsf{Event})} & \neg \mathsf{K} \ \phi \subseteq \mathsf{K} \ (\neg \mathsf{K} \ \phi) \end{split}$$

Though not explicitly included in their definition, it can be shown that knowledge operators retain knowledge generalisation and the Axiom K property. In each case, this follows from the preservation of semantic entailment by choosing an appropriate event family.

**Theorem 3.4.3.** Knowledge operators satisfy knowledge generalisation,  $\forall \phi \rightarrow \forall \mathsf{K} \phi$ .

*Proof.* Unfolding  $\forall \phi \to \forall \mathsf{K} \phi$  according to Definition 3.3.3, we are to prove the statement  $(\Pi_{(w:\mathsf{State})} \phi w) \to (\Pi_{(v:\mathsf{State})} \mathsf{K} \phi v)$ . We therefore proceed by assuming that  $\phi$  holds in any state, and continue to show  $\mathsf{K} \phi v$  for an arbitrary state v.

This is equivalent to the nullary entailment  $\mathsf{K} \circ \mathsf{elim}_0 \vDash \mathsf{K} \phi$ , from an event family indexed by the empty type,  $\mathbb{O} \to \mathsf{Event}$ . Applying preservation of semantic entailment, we have to show that  $\mathsf{elim}_0 \vDash \phi$ , whose conclusion is proved by our original assumption that  $\phi$  holds universally.

**Theorem 3.4.4.** *Knowledge operators satisfy Axiom K*,  $\mathsf{K} \ (\phi \Rightarrow \psi) \Rightarrow \mathsf{K} \ \phi \Rightarrow \mathsf{K} \ \psi$ .

*Proof.* The theorem statement  $\mathsf{K} \ (\phi \Rightarrow \psi) \subseteq \mathsf{K} \ \phi \Rightarrow \mathsf{K} \ \psi$  is judgementally equal to  $\Pi_{(w:\mathsf{State})} \ \mathsf{K} \ (\phi \Rightarrow \psi) \ w \to \mathsf{K} \ \phi \ w \to \mathsf{K} \ \psi \ w$ . We assume hypotheses  $\mathsf{K} \ (\phi \Rightarrow \psi) \ w$  and  $\mathsf{K} \ \phi \ w$  towards the goal  $\mathsf{K} \ \psi \ w$ .

Define the event family  $B = \operatorname{elim}_+ (\operatorname{elim}_1 (\phi \Rightarrow \psi)) (\operatorname{elim}_1 \phi) : 2 \to \operatorname{Event}$ , indexed by the Booleans. Applying preservation of semantic entailment with B, it suffices to show  $B \models \psi$  and  $\Pi_{(b:2)} \mathsf{K} (B \ b) w$ , the latter of which is simply a pairing of our hypotheses.

To prove  $B \vDash \psi$ , we may assume a dependent function of type  $\Pi_{(b:2)} B b w$ and use it to obtain a value of type  $\psi w$ . Applying the function to true and false, we have  $\phi w \rightarrow \psi w$  and  $\phi w$ , respectively, yielding  $\psi w$  by modus ponens.  $\Box$ 

### 3.5 Correspondence with Relational Semantics

To prove the knowledge operator semantics correspond to the relational Kripke semantics of epistemic logic, we define transformation functions between their types. First, the type of Kripke frames is defined in a similar way to the knowledge operator semantics, as a dependent tuple type collecting an equivalence relation and its properties.

**Definition 3.5.1.** An equivalence relation over some type A is a binary type family  $R: A \to A \to \mathsf{Type}$  satisfying the properties of reflexivity, symmetry, and transitivity.

A Kripke frame is an equivalence relation over the type of possible worlds, Kripke = EquivRel State.

We will define the transformation from a function K: Event  $\rightarrow$  Event to a relation State  $\rightarrow$  State  $\rightarrow$  Type without considering the knowledge operator and equivalence relation properties. Separately, we will demonstrate that if K is a knowledge operator, then the transformation must endow the equivalence relation properties on its output.

**Definition 3.5.2.** Transforming an event operator K, define the binary relation which holds between states whose events are not distinguished by K. That is, w R v when, K  $\phi w$  and K  $\phi v$  are propositionally equivalent for all  $\phi$ .

If classical reasoning is admitted, only one direction of the equivalence is sufficient, the other being derivable. For simplicity, this is the version of the transformation which will be used going forward.

$$\begin{split} R_{[\_]} : (\mathsf{Event} \to \mathsf{Event}) \to \mathsf{State} \to \mathsf{State} \to \mathsf{Type} \\ w \; R_{[\mathsf{K}]} \; v = \Pi_{(\phi:\mathsf{Event})} \; \mathsf{K} \; \phi \; w \to \mathsf{K} \; \phi \; v \end{split}$$

Note that, since this transformation is quantified over events, the resulting type universe must be at a higher level than the one **Event** itself resides in.

#### **Theorem 3.5.3.** If K is a knowledge operator, then $R_{[K]}$ is an equivalence relation.

*Proof.* Reflexivity and transitivity are immediate from fact that the underlying implication relation is a preorder — see Theorems 2.3.1 and 2.3.2.

To show symmetry, assume  $w R_{[K]} v$ , which is  $\Pi_{(\phi:\text{Event})} \mathsf{K} \phi w \to \mathsf{K} \phi v$  according to Definition 3.5.2. We are to prove the opposite direction for this implication, so also assume  $\mathsf{K} \phi v$  and aim to show  $\mathsf{K} \phi w$ . Using the classical principle of double-negation elimination, assume  $\mathsf{K} \phi w \to \mathbb{O}$ , towards a contradiction.

By Definition 3.3.2, this assumption is equal to  $(\neg \mathsf{K} \ \phi) \ w$ . Since  $\mathsf{K}$  is a knowledge operator, using the Axiom 5 property we can derive  $\mathsf{K} (\neg \mathsf{K} \ \phi) \ w$ , from which we can conclude  $\mathsf{K} (\neg \mathsf{K} \ \phi) \ v$  by instantiating our first hypothesis at index  $\neg \mathsf{K} \ \phi$ . By Axiom T, this implies  $(\neg \mathsf{K} \ \phi) \ v$ , which can be written  $\mathsf{K} \ \phi \ v \to \mathbb{O}$ , but this contradicts our assumption of  $\mathsf{K} \ \phi \ v$ . Therefore,  $\mathsf{K} \ \phi \ w$  is proved by contradiction.

We now come to defining the opposite transformation, from a binary relation  $R: \text{State} \rightarrow \text{State} \rightarrow \text{Type}$  into an operator  $\text{Event} \rightarrow \text{Event}$ .

**Definition 3.5.4.** Transforming a binary relation R, define the event operator whose knowledge is preserved by R. That is,  $\mathsf{K} \phi$  holds in state w when  $\phi$  holds in all states that are related to w.

 $\mathsf{K}_{[..]} : (\mathsf{State} \to \mathsf{State} \to \mathsf{Type}) \to \mathsf{Event} \to \mathsf{Event}$  $\mathsf{K}_{[R]} \ \phi \ w = \Pi_{(v:\mathsf{State})} \ w \ R \ v \to \phi \ v$ 

We demonstrate each knowledge operator property of  $\mathsf{K}_{[R]}$  only assuming the relevant equivalence relation properties in each case.

**Lemma 3.5.5.** Without any constraints on the relation,  $K_{[R]}$  preserves semantic entailment.

*Proof.* Assume for some event family  $\Phi : A \to \mathsf{Event}$  and event  $\psi$  that  $\Phi \vDash \psi$ . We want to reach the conclusion  $\mathsf{K}_{[R]} \circ \Phi \vDash \mathsf{K}_{[R]} \psi$ , which, applying Definitions 3.4.1

and 3.5.4, is the following statement.

$$\Pi_{(w:\mathsf{State})} (\Pi_{(a:A)} \Pi_{(v:\mathsf{State})} w R v \to \Phi a v) \to \Pi_{(v:\mathsf{State})} w R v \to \psi v$$

With additional hypotheses  $(\prod_{(a:A)} \prod_{(v:State)} w \ R \ v \to \Phi \ a \ v)$  and  $w \ R \ v$ , we need to show  $\psi \ v$ . This follows from the original assumption of  $\Phi \vDash \psi$  if we can show  $\Phi \ a \ v$ . But this is proved by our first additional hypothesis, instantiated at indices a and v, as we have already assumed  $w \ R \ v$ .

**Lemma 3.5.6.** When R is reflexive,  $K_{[R]}$  satisfies Axiom T.

*Proof.* Assume  $\mathsf{K}_{[R]} \phi w$  with the goal of proving  $\phi w$ . By Definition 3.5.4, the assumption expands to  $\Pi_{(v:\mathsf{State})} w R v \to \phi v$ , so instantiate it at w to obtain  $w R w \to \phi w$ . Our desired  $\phi w$  therefore follows from the reflexivity of R.  $\Box$ 

**Lemma 3.5.7.** When R is transitive,  $K_{[R]}$  satisfies Axiom 4.

*Proof.* Assume  $\mathsf{K}_{[R]} \phi w$ , which, as above, is  $\Pi_{(v:\mathsf{State})} w R v \to \phi v$ . We want to reach the conclusion  $\mathsf{K}_{[R]} (\mathsf{K}_{[R]} \phi) w$ , which is judgementally equal to the following.

$$\Pi_{(v:\mathsf{State})} w R v \to \Pi_{(u:\mathsf{State})} v R u \to \phi u$$

Assume  $w \ R \ v$  and  $v \ R \ u$  and use transitivity to obtain  $w \ R \ u$ . The result  $\phi \ u$  is then proved by applying the original assumption at index u.

**Lemma 3.5.8.** When R is symmetric and transitive,  $K_{[R]}$  satisfies Axiom 5.

*Proof.* Assume  $\neg \mathsf{K}_{[R]} \phi w$ . After unfolding  $\mathsf{K}_{[R]} (\neg \mathsf{K}_{[R]} \phi) w$ , we are to derive a contradiction.

$$\Pi_{(v:\mathsf{State})} \ w \ R \ v \to (\Pi_{(u:\mathsf{State})} \ v \ R \ u \to \phi \ u) \to \mathbb{O}$$

If we can prove  $\mathsf{K}_{[R]} \phi w$ , which is judgementally equal to  $\Pi_{(t:\mathsf{State})} w R t \to \phi t$ , this will contradict our assumption. To this end, assume additional hypotheses w R v, w R t, and  $\Pi_{(u:\mathsf{State})} v R u \to \phi u$ , setting out to prove  $\phi t$ .

Using symmetry and transitivity with the first two of these, we derive v R t. Then, applying the third hypothesis at index t reaches the desired result.

**Theorem 3.5.9.** When R is an equivalence relation,  $K_{[R]}$  is a knowledge operator.

*Proof.* Aggregating Lemmas 3.5.5 to 3.5.8 gives the full theorem.  $\Box$ 

Theorems 3.5.3 and 3.5.9 show that transforming knowledge operators and equivalence relations always yields equivalence relations and knowledge operators, respectively. For a complete correspondence between the two semantics, we furthermore want to show that the transformations are inverse, that knowledge operators and equivalence relations are isomorphic.

The isomorphism will be established up to propositional equivalence. For example, we won't show  $\mathsf{K} \equiv \mathsf{K}_{[R_{[\mathsf{K}]}]}$ , using propositional equality to state they are actually identical objects. Rather, we will show that the types  $\mathsf{K} \phi w$  and  $\mathsf{K}_{[R_{[\mathsf{K}]}]} \phi w$  are coinhabited — if either can be proved, then so can the other.

**Theorem 3.5.10.** When  $R_{[]}$  is applied to a knowledge operator,  $K_{[]}$  is its inverse, up to propositional equivalence. That is,  $K \phi \approx K_{[R_{[K]}]} \phi$ .

*Proof.* The left-to-right direction is  $\mathsf{K} \ \phi \subseteq \mathsf{K}_{[R_{[\mathsf{K}]}]} \psi$ , so assume  $\mathsf{K} \ \phi \ w$ . The conclusion is judgementally equal to the following statement.

$$\Pi_{(v:\mathsf{State})} \ (\Pi_{(\psi:\mathsf{Event})} \ \mathsf{K} \ \psi \ w \to \mathsf{K} \ \psi \ v) \to \phi \ v$$

With the additional assumption of  $\Pi_{(\psi:\mathsf{Event})} \mathsf{K} \psi w \to \mathsf{K} \psi v$  instantiated at  $\phi$ , we can derive  $\mathsf{K} \phi v$ , and then  $\phi v$  follows from the Axiom T property.

For the right-to-left direction, assume  $\mathsf{K}_{[R_{[\mathsf{K}]}]} \phi w$  — this is the same as the conclusion for the left-to-right direction, so its full definition is above. To show  $\mathsf{K} \phi w$ , we will use preservation of semantic entailment using as event family all those events which are known in w. Formally, we define the event family  $\Phi = \lambda(p : \Sigma_{(\psi:\mathsf{Event})} \mathsf{K} \psi w)$ .  $\mathsf{K} (\mathsf{elim}_{\Sigma} (\lambda \psi, \lambda k, \psi))$ . It is indexed by dependent pairs,  $\langle \psi, k \rangle$ , where k is a proof of  $\mathsf{K} \psi w$ , with the output being the event  $\mathsf{K} \psi$ .

To see that our premise  $\mathsf{K}_{[R_{[\mathsf{K}]}]} \phi w$  implies the entailment  $\Phi \vDash \phi$ , assume  $\Pi_{(p:\Sigma_{(\psi:\mathsf{Event})} \mathsf{K} \psi w)} \Phi p v$ . Then, given a proof  $k : \mathsf{K} \psi w$ , we can instantiate the assumption at index  $\langle \psi, k \rangle$ , yielding  $\Phi k v$ , which is  $\mathsf{K} \psi v$ . Applying the premise at event  $\phi$  therefore proves  $\phi v$ .

Finally, we must show  $\Pi_{(p:\Sigma_{(\psi:\mathsf{Event})} \mathsf{K} \psi w)} (\mathsf{K} \circ \Phi) p w$ . Given an arbitrary index  $\langle \psi, k \rangle$  where  $k : \mathsf{K} \psi w$ , this is  $\mathsf{K} (\mathsf{K} \psi) w$ , so we apply the Axiom 4 property to the second element of the pair.

**Theorem 3.5.11.** When  $\mathsf{K}_{[]}$  is applied to an equivalence relation,  $R_{[]}$  is its inverse, up to propositional equivalence. That is,  $w \ R \ v \leftrightarrow w \ R_{[\mathsf{K}_{[R]}]} \ v$ .

*Proof.* In the left-to-right direction, expanding only the outer transformation, we want to demonstrate that the knowledge information of  $K_{[R]}$  can be transported between states related by R.

$$w \mathrel{R} v \to \Pi_{(\phi:\mathsf{Event})} \mathrel{\mathsf{K}}_{[R]} \phi \mathrel{w} \to \mathrel{\mathsf{K}}_{[R]} \phi \mathrel{v}$$

But, in fact, this is an instance the Axiom 4 property specialised to  $K_{[R]}$ , whose use is justified by Theorem 3.5.9. This can be seen by selectively expanding

definitions in the statement  $\mathsf{K}_{[R]} \phi \subseteq \mathsf{K}_{[R]} (\mathsf{K}_{[R]} \phi)$ .

$$\Pi_{(w:\mathsf{State})} \mathsf{K}_{[R]} \phi w \to \Pi_{(v:\mathsf{State})} w R v \to \mathsf{K}_{[R]} \phi v$$

For the right-to-left direction, assume  $w R_{[K_{[R]}]} v$ . After fully-unfolding, this is judgementally equal to the following statement.

$$\Pi_{(\phi:\mathsf{Event})} \ (\Pi_{(u:\mathsf{State})} \ w \ R \ u \to \phi \ u) \to \Pi_{(u:\mathsf{State})} \ v \ R \ u \to \phi \ u$$

Instantiating with event  $\phi = \lambda t$ . w R t and state v, its conclusion becomes our desired w R v.

$$(\Pi_{(u:\mathsf{State})} \ w \ R \ u \to w \ R \ u) \to v \ R \ v \to w \ R \ v$$

We need only provide the identity function  $\mathsf{id}_{(wRu)}$  and a proof of the reflexive property,  $v \ R \ v$ .

This isomorphism justifies working with either representation of knowledge equivalently, whichever is more convenient at each point. Given a knowledge operator or an equivalence relation over states, we can leverage our transformations to obtain the other, complete with all of its expected properties, and preserving the truth value when transforming back.

In Chapter IV, we shall work with the shallow embedding defined in this chapter to define the concept of *common knowledge*. A type of agents will be assumed, with an equivalence relation postulated for each. The agents' knowledge operators are defined directly using the  $K_{[\_]}$  transformation, and the correspondence between the two semantics will be essential.

# Chapter IV

# **Coinductive Common Knowledge**

In this chapter, we introduce the concept of *common knowledge* in epistemic modal logic with multiple agents, and show how the relational semantics is extended to support the new operator. We will then use the embedding defined in Chapter III to define a common knowledge operator directly, using the feature of coinductive types in type theory.

This coinductive common knowledge operator will be shown to correspond to the intuitive notion of iterated universal knowledge. Furthermore, the common knowledge operator induced by the relation semantics will also be defined in the embedding, and proved equivalent to the coinductive definition.

An Agda formalisation of both definitions for common knowledge and their equivalence is included in Appendix A.3

## 4.1 Background

Common knowledge is a concept in epistemic logic when it is extended to more than one agent. If a fact  $\phi$  is common knowledge among the agents, it not only implies that it is known universally, but also that the agents can reason about each others' knowledge of  $\phi$  to an arbitrary degree. While common knowledge was first given a mathematical definition by Robert Aumann [2], the notion was discussed much earlier in the philosophical literature [41].

Assumptions of common knowledge permeate society. For example, it can be assumed that it is common knowledge among UK drivers that they are to drive on the left side of the road. Not only do drivers assume that all other drivers know this, they also assume that all drivers assume this of all other drivers, and so on. If it were announced that the nation would suddenly switch to driving on the right, common knowledge of this fact would have to be established before drivers could feel reasonably safe on the road again. Formally, in the syntax of epistemic logic, given a set of agents A, there is a modal operator  $\mathsf{K}_a$  for each agent  $a \in A$ , and we have new modal operators for universal knowledge,  $\mathsf{EK}$  read "everyone knows", and for common knowledge,  $\mathsf{CK}$ . Universal knowledge of  $\phi$  is defined as the conjunction of the knowledge for all agents,  $\mathsf{EK} \phi = \bigwedge_{a \in A} \mathsf{K}_a \phi$ , while common knowledge of  $\phi$  is understood as an infinite conjunction,  $\mathsf{CK} \phi = \bigwedge_{n \in \mathbb{N}} \mathsf{EK}^n \phi$ , combining all finite iterations of  $\mathsf{EK}$  [4]. That is, from an assumption of  $\mathsf{CK} \phi$ , we can derive  $\mathsf{EK} \phi$ , and therefore  $\mathsf{K}_a \phi$  for any agent a, as well as  $\mathsf{EK} (\mathsf{EK} \phi)$ ,  $\mathsf{EK} (\mathsf{EK} (\mathsf{EK} \phi))$ , and so on ad infinitum.

Let us explore the semantics of common knowledge in the context of a famous logic puzzle, known as the *Blue Eyes Puzzle* or the *Muddy Children Puzzle*, among other variants. The following is a rephrasing of the first version I heard, from Randall Munroe's XKCD website [46], but the puzzle was widely shared as at least as early as the 1950s [10, 26].

A group of logicians are performing an experiment on a remote island. Every night, a ferry arrives which will allow them to leave, but only those who can accurately determine the colour of their eyes are be allowed to board. While every logician does want to leave, they are committed to the experiment and will never guess their eye colour, only leaving when they can be certain by logical deduction. Due to the experimental parameters, the logicians cannot communicate in any way, but each is able to see the others, and each trusts in the deduction ability of the others.

One day, impatient after no logician has made an attempt to leave, the experiment organiser makes an announcement, "at least one of you has blue eyes". In fact, there are three total blue-eyed logicians, and everyone already knows the fact that was announced — those without blue eyes can see the three with blue eyes, and each of those three can see two others. Nevertheless, all three blue-eyed logicians leave on the third ferry after the announcement. Why?

The puzzle involves common knowledge in two ways. First, we may assume that the setup of the puzzle is common knowledge among all the logicians. For example, each will leave the island when deducing their own eye colour, but they also assume this of the others, and assume that this is common knowledge. This is similar to Aumann's assumption of common knowledge of *rationality* in game theory [3] — each player will maximise their own payoff, assumes the other players will do the same, and that this is common knowledge.

Second, when the announcement is made, it might initially seem that no additional information is given since the statement was already known universally. However, assuming the announcement is trusted, and the trust is itself common knowledge, what was universal knowledge before the announcement becomes common knowledge after the announcement.

Let  $\phi$  be the statement of the announcement, that there is at least one blueeyed logician. Since each logician can see at least two others with blue eyes, and the rationality of the logicians is common knowledge, they are all able to deduce EK  $\phi$  before the announcement. However, after the announcement, CK  $\phi$  holds and only then are they able to make the crucial deduction, EK (EK  $\phi$ ). This can be generalised by the following inductive argument.

**Theorem 4.1.1.** When there are  $n \ge 1$  blue-eyed logicians, they will all leave on the  $n^{th}$  ferry after the announcement.

*Proof.* If there is only one blue-eyed logician, before the announcement he doesn't know  $\phi$ . After the announcement,  $\mathsf{EK} \phi$  holds, so the blue-eyed logician can deduce  $\phi$ . Seeing no one else with blue eyes, the statement must refer to him, so he leaves on the first ferry.

If there are two blue-eyed logicians, they each know  $\phi$  before the announcement, but not  $\mathsf{EK} \phi$  — not knowing they have blue eyes themselves, they each reason that they could be in a possible world corresponding to one blue-eyed case. After the announcement,  $\mathsf{EK} (\mathsf{EK} \phi)$  holds, so each can deduce  $\mathsf{EK} \phi$ . In particular, they each now know that  $\phi$  is known by the other. The next day, upon seeing the other still present, they rule out that possible world, leaving only the one where they both have blue eyes, so they both leave on the second ferry.

In general, when there are  $n \geq 2$  blue-eyed logicians, each of them knows  $\mathsf{EK}^{n-2} \phi$  from their own observations, but  $\mathsf{EK}^{n-1} \phi$  isn't universally known. After the announcement,  $\mathsf{EK}^n \phi$  holds, so they can now deduce  $\mathsf{EK}^{n-1} \phi$ . Given this, all blue-eyed logicians reason by inductive hypothesis that if there are only n-1 of them, they will all leave on night n-1. However, n-1 nights pass without anyone leaving, falsifying the premise of the hypothesis — there must actually be at least n blue-eyed logicians. Only seeing n-1, they each conclude there must be exactly n, which must include them, therefore all n leave the next night, on the  $n^{\text{th}}$  ferry.

Recall that in the semantics of epistemic logic, the knowledge of an agent may be modelled by an equivalence relation over possible worlds, called a Kripke frame. In multi-agent scenarios, each agent has their own relation independent of the others. We can describe the knowledge content of the puzzle scenario by giving a knowledge relation for each logician.

The relevant information about the state of the world is whether each logician has blue eyes or not. We can represent this information as a tuple of Booleans, for example  $\langle \top, \top, \top \rangle$  for three total logicians, each of whom have blue eyes. In this state, since the first logician doesn't know his own eye colour, he could be in state  $\langle \bot, \top, \top \rangle$  as far as he is aware, so these states are related in the model, but he can distinguish states which differ by the other components, so these are not related.

In general, each logician's knowledge relation corresponds to an axis of an n-dimensional cube [24]. For the three-logician case, this is can be shown by the following three graphs, leaving the reflexive connections implicit.

Common knowledge is itself interpreted with its own knowledge relation. An event which is common knowledge must be known universally, so states must be related by the common knowledge relation when they are related by any one of the individual agents' knowledge relations. Set-theoretically, this is the union of all agents' relations, intuitively taking the union of their ignorance about the world. For the above three-logician example, this operation completes the cube.

$$\begin{array}{c|c} \langle \bot, \top, \top \rangle & \longrightarrow \langle \top, \top, \top \rangle \\ \hline \\ \langle \bot, \top, \bot \rangle & + & \langle \top, \top, \bot \rangle \\ & & & & \\ & & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & & \\ & & & \\ & & & \\ & &$$

This is the relational interpretation of the universal knowledge operator, EK. Observe that states with two or more blue-eyed logicians are not directly connected to  $\langle \perp, \perp, \perp \rangle$ , which corresponds to the fact that everyone already knows there is at least one blue-eyed logician in these states, as is stated in the puzzle text. The relation is reflexive and symmetric, but not transitive, which implies that the corresponding operator lacks the introspective properties corresponding to Axioms 4 and 5 of the S5 modal logic. From only an assumption of EK  $\phi$ , we cannot in general conclude EK (EK  $\phi$ ).

To get a full accounting of common knowledge, which includes the Axiom 4 and 5 properties, we further take the transitive closure. In the example, the

relation will be strongly connected, all states in a single equivalence class. It is as if there is a hypothetical, maximally-ignorant agent whose knowledge relation is the common knowledge relation, unable to observe a difference between any pair of states.

However, following the announcement which makes state  $\langle \perp, \perp, \perp \rangle$  impossible, the connections from this state in the logicians' individual relations are severed. In states where there is a single blue-eyed logician, that logician is now able to precisely determine the state, while the other two logicians still consider another state possible.

Constructing the updated common knowledge relation by union and transitive closure, state  $\langle \perp, \perp, \perp \rangle$  now forms its own equivalence class. This indicates that it has become common knowledge that there is at least one blue-eyed logician in the other states. Even our hypothetical agent can now distinguish this state from the others. In the following graph for this relation, both the reflexive and transitive connections are left implicit.



The first night after the announcement, world states with only one blue-eyed logician become impossible from the perspective of logicians in states with at least two. In these states, it has therefore become common knowledge that there are at least two blue-eyed logicians. The common knowledge relation now has singleton equivalence classes for all states with fewer than two blue-eyed logicians.

$$\begin{array}{c|c} \langle \bot, \top, \top \rangle & \longrightarrow & \langle \top, \top, \top \rangle \\ \langle \bot, \top, \bot \rangle & \langle \top, \top, \bot \rangle & & \\ \langle \bot, \bot, \top \rangle & & \langle \top, \bot, \top \rangle \\ \langle \bot, \bot, \bot \rangle & & \langle \top, \bot, \bot \rangle \end{array}$$

This pattern continues after the second night, with the entire world state becoming common knowledge in the three-logician case, each state in its own singleton equivalence class. In the general case with arbitrarily many logicians, n of which have blue eyes, common knowledge of the value of n is established n-1 nights after the announcement.

## 4.2 Coinductive Definition

We now return to the embedding of epistemic logic into type theory defined in Chapter III. To define a common knowledge operator in the embedding, we need to add type of agents. We postulate this type abstractly, making only the assumption that it is non-empty. That is, we assume an element  $\langle \text{Agent, agent} \rangle$ of type  $\Sigma_{(\text{Agent:Type})}$  Agent. We therefore have a type, Agent, which has at least one concrete element, agent : Agent.

The isomorphism between the knowledge operator semantics and the relational semantics was established, so we may equivalently work with either representation. Here we assume equivalence relations as primitive, each agent aequipped with  $\propto_a$ : EquivRel State.

**Definition 4.2.1.** The agents' knowledge operators are defined by the isomorphism, with Theorem 3.5.9 validating they have the appropriate properties.

$$\mathsf{K}: \mathsf{Agent} o \mathsf{Event} o \mathsf{Event}$$
  
 $\mathsf{K}_a = \mathsf{K}_{[\propto_a]}$ 

With a type of agents now equipped with knowledge operators, universal knowledge among them has an obvious definition.

**Definition 4.2.2.** An event  $\phi$  is universal knowledge when "everyone knows"  $\phi$ . That is, each agent knows  $\phi$  individually.

 $\begin{array}{l} \mathsf{EK}:\mathsf{Event}\to\mathsf{Event}\\ \mathsf{EK}\;\phi\;w=\Pi_{(a:\mathsf{Agent})}\;\mathsf{K}_a\;\phi\;w \end{array}$ 

The EK operator is not a full knowledge operator. It preserves semantic entailments and satisfies the Axiom T property, but not the two introspective properties corresponding to Axioms 4 and 5. Intuitively, the relational representation of the operator is the union of all agents' individual relations — it is reflexive and symmetric, but not necessarily transitive.

**Lemma 4.2.3.** EK preserves semantic entailment,  $\Phi \vDash \psi \rightarrow \mathsf{EK} \circ \Phi \vDash \mathsf{EK} \psi$ .

*Proof.* Assume an event family  $\Phi : T \to \mathsf{Event}$  such that  $\Phi \vDash \psi$ . Given hypothesis  $\Pi_{(t:T)} \mathsf{EK} (\Phi t) w$ , we need to show  $\mathsf{EK} \psi w$ . By Definitions 4.2.1 and 4.2.2, this

is the following statement.

$$\Pi_{(a: \mathsf{Agent})} \Pi_{(v: \mathsf{State})} w \propto_a v \to \psi v$$

So we assume  $w \propto_a v$  and apply the entailment  $\Phi \vDash \psi$  at index v, the goal becoming  $\Phi t v$ . From our earlier hypothesis, we know EK  $(\Phi t) w$ , so can derive  $\mathsf{K}_a \ (\Phi t) w$  by instantiating universal knowledge specialised to agent a.

Using Axiom 4 as a transportation principle as in Theorem 3.5.11, we can transport this knowledge to state v as these states are bridged by  $w \propto_a v$ . Since  $\mathsf{K}_a \ (\Phi \ t) \ v$ , Axiom T validates  $\Phi \ t \ v$  — if a knows  $\Phi \ t$  at state v, the event must actually occur there.

**Lemma 4.2.4.** EK satisfies Axiom T, EK  $\phi \subseteq \phi$ .

*Proof.* Given  $\mathsf{EK} \phi w$ , instantiate this at the agent constant to derive  $\mathsf{K}_{\mathsf{agent}} \phi w$ . The conclusion  $\phi w$  then follows from the Axiom T property of agent's knowledge operator.

With universal knowledge defined, we can give common knowledge a definition as an operator on events. The intuitive idea of common knowledge as the infinite conjunction  $\mathsf{EK} \ \phi \land \mathsf{EK} \ (\mathsf{EK} \ \phi) \land \ldots$  can be realised directly using a coinductive type.

**Definition 4.2.5.** Common knowledge of an event  $\phi$  implies that it is universal knowledge and, corecursively, common knowledge of the universal knowledge.

 $\begin{array}{l} \mbox{coinductive CK}: \mbox{Event} \to \mbox{Event} \\ \mbox{introCK}: \Pi_{(\phi: \mbox{Event})} \ \mbox{EK} \ \phi \wedge \mbox{CK} \ (\mbox{EK} \ \phi) \subseteq \mbox{CK} \ \phi \end{array}$ 

Expanding the type of the constructor using Definitions 3.3.2 and 3.3.3, it can be seen that the canonical way to prove CK  $\phi$  w is to provide proofs of EK  $\phi$  w and CK (EK  $\phi$ ) w collected as a pair.

introCK :  $\Pi_{(\phi:\mathsf{Event})} \Pi_{(w:\mathsf{State})} \mathsf{EK} \phi w \times \mathsf{CK} (\mathsf{EK} \phi) w \to \mathsf{CK} \phi w$ 

This coinductive definition resembles that of infinite streams. They each have only one constructor with a single recursive branch, carrying some additional data at each node. However, the type of the data carried by streams is the same at every node, where the common knowledge constructor is a dependent function which changes the type at each node, each step having one more application of EK than the previous.

Since a proof of  $\mathsf{CK} \phi$  is constructed from a conjunction, it is immediate that either conjunct may be projected from the underlying product type. Corollary 4.2.6. CK  $\phi \subseteq$  EK  $\phi$ 

#### Corollary 4.2.7. CK $\phi \subseteq$ CK (EK $\phi$ )

A similar property to Corollary 4.2.7, swapping the operators in the conclusion, states common knowledge of an event implies that it's universally known to be common knowledge. However, this is not immediate from the definition, requiring a coinductive proof.

#### Lemma 4.2.8. CK $\phi \subseteq$ EK (CK $\phi$ )

*Proof.* In a coinductive proof of common knowledge, we may recursively use the statement we are proving as long as its use is guarded by an application of constructor introCK to prevent circularity. We therefore have the coinductive hypothesis  $\Pi_{(\phi:\text{Event})}$  CK  $\phi \subseteq$  EK (CK  $\phi$ ).

We assume  $\mathsf{CK} \phi w$  and, from Definitions 4.2.1 and 4.2.2,  $w \propto_a v$ . From the first assumption and Corollary 4.2.7, we obtain  $\mathsf{CK} (\mathsf{EK} \phi) w$ . The conclusion is the statement  $\mathsf{CK} \phi v$ , so apply intro $\mathsf{CK}$ , thus requiring proofs for  $\mathsf{EK} \phi v$  and  $\mathsf{CK} (\mathsf{EK} \phi) v$ .

To show EK  $\phi v$  from CK (EK  $\phi$ ) w, use Corollary 4.2.6 to reach EK (EK  $\phi$ ) wand therefore its specialisation, K<sub>a</sub> (EK  $\phi$ ) w. By the assumption of  $w \propto_a v$ , we use Axiom 4 to transport the knowledge of EK  $\phi$  from state w to v.

To show CK (EK  $\phi$ ) v from CK (EK  $\phi$ ) w, guardedness is satisfied, so we may use the coinductive hypothesis at index EK  $\phi$  to obtain EK (CK (EK  $\phi$ )) w and its specialisation, K<sub>a</sub> (CK (EK  $\phi$ )) w. By the same reasoning as above, the known event may then be transported to state v.

By repeatedly applying Corollary 4.2.7 followed by an application of Corollary 4.2.6, we can show that EK can be iterated on  $\phi$  any number of times from an assumption of CK  $\phi$ , which matches the intuition of common knowledge as iterated universal knowledge. To demonstrate this formally, we define the event family recCK  $\phi$ , indexed by the natural numbers, which contains all those events defined by iterating EK — recCK  $\phi$  0 =  $\phi$ , recCK  $\phi$  1 = EK  $\phi$ , recCK  $\phi$  2 = EK (EK  $\phi$ ), and so on.

**Definition 4.2.9.** The event recCK  $\phi$  *n* is defined by recursively iterating the function EK *n* times over  $\phi$ .

$$\begin{split} \mathsf{recCK} &: \mathsf{Event} \to \mathbb{N} \to \mathsf{Event} \\ \mathsf{recCK} \ \phi \ \mathsf{zero} &= \phi \\ \mathsf{recCK} \ \phi \ (\mathsf{succ} \ n) = \mathsf{recCK} \ (\mathsf{EK} \ \phi) \ n \end{split}$$
#### **Theorem 4.2.10.** The event CK $\phi$ semantically entails event family recCK $\phi$ .

*Proof.* Proceed by induction on the indices of event family recCK  $\phi$ . The base case is  $\Pi_{(\phi:\mathsf{Event})} \mathsf{CK} \phi \subseteq \mathsf{recCK} \phi 0$ , whose conclusion reduces to  $\phi$ . Assume CK  $\phi$ , applying Corollary 4.2.6 to obtain EK  $\phi$ , followed by the Axiom T property of EK, Lemma 4.2.4, to conclude  $\phi$ .

We have inductive hypothesis  $\Pi_{(\phi:\mathsf{Event})} \mathsf{CK} \phi \subseteq \mathsf{recCK} \phi n$  for the inductive case. Assume  $\mathsf{CK} \phi$  and apply Corollary 4.2.7 to obtain  $\mathsf{CK} (\mathsf{EK} \phi)$ . Use the inductive hypothesis at event  $\mathsf{EK} \phi$  to conclude  $\mathsf{recCK} (\mathsf{EK} \phi) n$ , which is equal to  $\mathsf{recCK} \phi$  (succ n).

We can also show the reverse entailment as another example of a coinductive proof for common knowledge. Therefore the equivalence of event  $\mathsf{CK} \phi$  and event family recCK  $\phi$  is established.

#### **Theorem 4.2.11.** The event family recCK $\phi$ semantically entails event CK $\phi$ .

*Proof.* Assume, as the coinductive hypothesis,  $\Pi_{(\phi:\mathsf{Event})} \operatorname{\mathsf{recCK}} \phi \vDash \mathsf{CK} \phi$ . Then assume  $\Pi_{(n:\mathbb{N})} \operatorname{\mathsf{recCK}} \phi n$  and use intro $\mathsf{CK}$  to construct a term of coinductive type  $\mathsf{CK} \phi$ .

We need to provide a pair of proofs, the first of which,  $\mathsf{EK} \phi$ , is established by applying our assumption at index 1. For the second,  $\mathsf{CK} (\mathsf{EK} \phi)$ , guardedness is satisfied, so we apply the coinductive hypothesis at event  $\mathsf{EK} \phi$ . We then need to show  $\Pi_{(n:\mathbb{N})}$  recCK ( $\mathsf{EK} \phi$ ) n, which is the same as  $\Pi_{(n:\mathbb{N})}$  recCK  $\phi$  (succ n), yielded by applying the assumption at index succ n

Unlike the universal knowledge operator it's based on, the coinductive common knowledge operator is a true knowledge operator, satisfying all required properties. While each property may be proved directly, it is more convenient to leverage the equivalence with the relational definition of common knowledge. The knowledge operator proof will be shown with this equivalence in Section 4.4.

## 4.3 Relational Definition

The relational semantics of epistemic logic interprets common knowledge among a group of agents as the equivalence relation that is the transitive closure of the union of the individual agents' relations. This notion can be defined in our embedding directly as an inductive relation, denoted by  $\propto_{Agent}$ . The knowledge operator semantics then follows from Theorem 3.5.9 provided we show  $\propto_{Agent}$  is an equivalence relation. **Definition 4.3.1.** The statement  $w \propto_{\mathsf{Agent}} v$  may be derived when there is an agent *a* such that  $w \propto_a v$ . Alternatively, it may be derived transitively.

 $\begin{array}{l} \mathsf{inductive} \propto_{\mathsf{Agent}} : \mathsf{State} \to \mathsf{State} \to \mathsf{Type} \\ \\ \mathsf{union}_{\propto_{\mathsf{Agent}}} : \Pi_{(a:\mathsf{Agent})} \ \Pi_{(w,v:\mathsf{State})} \ w \propto_{a} v \to w \propto_{\mathsf{Agent}} v \\ \\ \\ \mathsf{trans}_{\propto_{\mathsf{Agent}}} : \Pi_{(w,v,u:\mathsf{State})} \ w \propto_{\mathsf{Agent}} v \to v \propto_{\mathsf{Agent}} u \to w \propto_{\mathsf{Agent}} u \end{array}$ 

**Definition 4.3.2.** The relational common knowledge operator, rCK, is yielded when the  $K_{[-]}$  transformation is applied to relation  $\propto_{Agent}$ .

 $\label{eq:rCK} \begin{array}{l} \mathsf{rCK}:\mathsf{Event}\to\mathsf{Event}\\ \\ \mathsf{rCK}=\mathsf{K}_{[\propto_{\mathsf{Agent}}]} \end{array}$ 

For rCK to be a knowledge operator, we must establish that  $\propto_{Agent}$  is indeed an equivalence relation. Transitivity is trivial, while reflexivity and symmetry remain to be proved.

**Lemma 4.3.3.** Reflexivity,  $w \propto_{\mathsf{Agent}} w$ .

*Proof.* Since there is at least one agent, witnessed by agent, whose knowledge relation is reflexive, we have  $w \propto_{\text{agent}} w$ . Applying the union<sub> $\propto_{\text{Agent}}$ </sub> constructor lifts this to  $w \propto_{\text{Agent}} w$ , as desired.

**Lemma 4.3.4.** Symmetry,  $w \propto_{\mathsf{Agent}} v \rightarrow v \propto_{\mathsf{Agent}} w$ .

*Proof.* Proceed by structural induction on the assumption of  $w \propto_{\mathsf{Agent}} v$ . We have two cases to consider, the base case when the assumption is constructed by  $\mathsf{union}_{\propto_{\mathsf{Agent}}}$ , and the inductive case when it is constructed by  $\mathsf{trans}_{\propto_{\mathsf{Agent}}}$ .

In the base case, there must be an agent a such that  $w \propto_a v$ , and whose knowledge relation is symmetric, so  $v \propto_a w$ . Applying  $\mathsf{union}_{\propto_{\mathsf{Agent}}}$  shows  $v \propto_{\mathsf{Agent}} w$ .

In the inductive case, we have assumptions  $w \propto_{\text{Agent}} u$  and  $u \propto_{\text{Agent}} v$ , each with an inductive hypothesis,  $u \propto_{\text{Agent}} w$  and  $v \propto_{\text{Agent}} u$ , respectively. Recombining the inductive hypotheses in the opposite order with the trans<sub> $\propto_{\text{Agent}}</sub> constructor therefore gives <math>v \propto_{\text{Agent}} w$ .</sub>

**Theorem 4.3.5.**  $\propto_{\text{Agent}}$  is an equivalence relation.

*Proof.* Transitivity is by definition of  $trans_{\propto Agent}$ , while reflexivity and symmetry are shown in Lemmas 4.3.3 and 4.3.4.

Corollary 4.3.6. rCK is a knowledge operator.

In addition to the knowledge operator properties, rCK satisfies the additional properties expected of an interpretation for common knowledge. These are useful lemmas, corresponding to properties already shown for CK, which describe how assumptions of rCK interact with EK.

#### Lemma 4.3.7. rCK $\phi \subseteq$ EK $\phi$

*Proof.* Assume rCK  $\phi$  w and  $w \propto_a v$ , setting out to prove  $\phi$  v. Using constructor union<sub> $\propto Agent$ </sub> with the second assumption, we obtain  $w \propto_{Agent} v$ . By Definition 4.3.2, the first assumption is equal to  $\Pi_{(v:State)} w \propto_{Agent} v \to \phi v$ , so applying this at index v reaches the desired conclusion.

#### Lemma 4.3.8. rCK $\phi \subseteq$ rCK (EK $\phi$ )

*Proof.* Assume rCK  $\phi$  w,  $w \propto_{\mathsf{Agent}} v$ , and  $v \propto_a u$ . Using constructor  $\mathsf{union}_{\propto_{\mathsf{Agent}}}$  with the last assumption, we obtain  $v \propto_{\mathsf{Agent}} u$ , and then, by  $\mathsf{trans}_{\propto_{\mathsf{Agent}}}$  with this and the second assumption,  $w \propto_{\mathsf{Agent}} u$ . Applying the first assumption, gives the desired result,  $\phi u$ .

#### Lemma 4.3.9. rCK $\phi \subseteq$ EK (rCK $\phi$ )

*Proof.* Assume rCK  $\phi$  w, w  $\propto_a v$ , and v  $\propto_{\text{Agent}} u$ . As above, use the  $\propto_{\text{Agent}}$  constructors with the second and third assumptions, then apply the first assumption to conclude  $\phi$  u.

#### 4.4 Coinductive and Relational Equivalence

We now come to proving the equivalence of the two operators we have defined for common knowledge,  $\mathsf{CK} \ \phi \approx \mathsf{rCK} \ \phi$ . There is one final lemma we need which provides a bridge between these two interpretations.

From Corollary 4.3.6, we know that the relational definition is a knowledge operator, in particular satisfying Axiom 4,  $rCK \phi \subseteq rCK (rCK \phi)$ . Recall that this property, for the  $K_a$  operator, implied a principle which allowed the agent's knowledge in one state to be transported to another, provided that these states were bridged by the agent's knowledge relation. Since the relational common knowledge operator is defined in the same way, using the  $K_{[\_]}$  transformation on an equivalence relation, we have a corresponding transportation principle for it.

$$\mathsf{rCK} \phi w \to w \propto_{\mathsf{Agent}} v \to \mathsf{rCK} \phi v$$

If we can show that coinductive common knowledge can also be transported across this relation, we have all we need to show the equivalence of the two operators. **Lemma 4.4.1.** Coinductive common knowledge can be transported across  $\propto_{\mathsf{Agent}}$ . That is,  $\mathsf{CK} \ \phi \ w \to w \propto_{\mathsf{Agent}} v \to \mathsf{CK} \ \phi \ v$ .

Proof. Proceed by structural induction on  $w \propto_{\mathsf{Agent}} v$ . In the base case,  $w \propto_{\mathsf{Agent}} v$  is constructed by  $\mathsf{union}_{\propto_{\mathsf{Agent}}}$ , so we need to show  $\mathsf{CK} \ \phi \ w \to w \propto_a v \to \mathsf{CK} \ \phi \ v$ . Assume  $\mathsf{CK} \ \phi \ w$  and  $w \propto_a v$ , then we will utilise Lemma 4.2.8 rather than directly proving  $\mathsf{CK} \ \phi \ v$  by coinduction. This tells us  $\mathsf{EK} \ (\mathsf{CK} \ \phi) \ w$ , which can be specialised  $\mathsf{K}_a \ (\mathsf{CK} \ \phi) \ w$ . Using Axiom 4 to transport this knowledge to state v, followed by Axiom T, we can conclude  $\mathsf{CK} \ \phi \ v$ .

In the inductive case,  $w \propto_{\mathsf{Agent}} v$  is constructed by  $\mathsf{trans}_{\propto_{\mathsf{Agent}}}$ , so we have assumptions  $\mathsf{CK} \phi w$ ,  $w \propto_{\mathsf{Agent}} u$ , and  $u \propto_{\mathsf{Agent}} v$ . The corresponding inductive hypotheses are  $\mathsf{CK} \phi w \to \mathsf{CK} \phi u$  and  $\mathsf{CK} \phi u \to \mathsf{CK} \phi v$ . We chain the inductive hypotheses to conclude  $\mathsf{CK} \phi v$ .

**Theorem 4.4.2.** Coinductive and relational common knowledge are equivalent, CK  $\phi \approx rCK \phi$ .

*Proof.* For the left-to-right direction, assume  $\mathsf{CK} \phi w$  and  $w \propto_{\mathsf{Agent}} v$ , towards the conclusion  $\phi v$ . By Lemma 4.4.1, we derive  $\mathsf{CK} \phi v$ . Since we have not yet proved  $\mathsf{CK}$  is a knowledge operator, we cannot apply Axiom T directly. Instead, we apply Corollary 4.2.6, to derive  $\mathsf{EK} \phi v$ , which has its own Axiom T principle, shown in Lemma 4.2.4.

For the right-to-left direction, assume rCK  $\phi$  w with the goal of showing CK  $\phi$  w. This direction is by coinduction, so apply the constructor introCK.

For the first field of the coinductive type, since we know rCK  $\phi w$ , Lemma 4.3.7 tells us EK  $\phi w$ . For the second field, use Lemma 4.3.8 to derive rCK (EK  $\phi$ ) w. We can then use the coinductive hypothesis indexed with event EK  $\phi$  to reach CK (EK  $\phi$ ) w, as desired.

We will finally show that CK is itself a knowledge operator. This follows from Theorem 4.4.2 and the status of rCK as a knowledge operator, Corollary 4.3.6.

Unfolding  $\mathsf{CK} \ \phi \approx \mathsf{rCK} \ \phi$  according to Definitions 3.3.2 and 3.3.3, we can see that it is a statement of propositional equivalence,  $\Pi_{(w:\mathsf{State})} \ \mathsf{CK} \ \phi \ w \leftrightarrow \mathsf{rCK} \ \phi \ w$ . If this were the stronger statement of propositional equality, we could appeal to Theorem 2.5.6 to substitute  $\mathsf{CK}$  by  $\mathsf{rCK}$  and the proof would be immediate, however propositional equivalence will only allow substitution between  $\mathsf{CK}$  and  $\mathsf{rCK}$  when one of these is the top-level connective of a hypothesis or the conclusion to be reached.

For example, to prove the Axiom 4 property,  $\mathsf{CK} \ \phi \subseteq \mathsf{CK} \ (\mathsf{CK} \ \phi)$ , after assuming the premise, we may substitute the conclusion  $\mathsf{CK} \ (\mathsf{CK} \ \phi) \ w$  by  $\mathsf{rCK} \ (\mathsf{CK} \ \phi) \ w$  using Theorem 4.4.2, but it will not allow direct substitution of the inner instance

of CK. To make further progress with the proof, we must make use of the definition and properties of rCK. We proceed in this fashion for each of the knowledge operator properties, making liberal use of Theorem 4.4.2 to be able to utilise the properties of rCK.

**Lemma 4.4.3.** CK preserves semantic entailment,  $\Phi \vDash \psi \rightarrow \mathsf{CK} \circ \Phi \vDash \mathsf{CK} \psi$ .

*Proof.* For an arbitrary event family  $\Phi : T \to \mathsf{Event}$ , assume hypotheses  $\Phi \vDash \psi$ and  $\Pi_{(t:T)} \mathsf{CK} (\Phi t) w$ . We are to prove  $\mathsf{CK} \psi w$ , but due to Theorem 4.4.2 we can equivalently prove  $\mathsf{rCK} \psi w$ . We may therefore assume  $w \propto_{\mathsf{Agent}} v$  with the goal becoming  $\psi v$ .

Applying the first hypothesis, it is sufficient to prove  $\Phi t v$  for some index t : T. We know  $\mathsf{CK} \ (\Phi t) w$ , which can be transported to state v using Lemma 4.4.1. Then, using the equivalence again, in the other direction, we have  $\mathsf{rCK} \ (\Phi t) v$ , and we may use Axiom T to reach the conclusion.

**Lemma 4.4.4.** CK satisfies Axiom T, CK  $\phi \subseteq \phi$ .

*Proof.* Assume  $\mathsf{CK} \phi w$  and use Theorem 4.4.2 to derive  $\mathsf{rCK} \phi w$ . The conclusion  $\phi w$  then follows from the Axiom T property of  $\mathsf{rCK}$ .

**Lemma 4.4.5.** CK satisfies Axiom 4, CK  $\phi \subseteq$  CK (CK  $\phi$ ).

*Proof.* Assume  $\mathsf{CK} \phi w$  and, using Theorem 4.4.2 with the conclusion as needed,  $w \propto_{\mathsf{Agent}} v$  and  $v \propto_{\mathsf{Agent}} u$ . To conclude  $\phi u$ , we chain Lemma 4.4.1 twice with our assumptions to reach  $\mathsf{CK} \phi u$ , and use the Axiom T property from Lemma 4.4.4.

**Lemma 4.4.6.** CK satisfies Axiom 5,  $\neg \mathsf{CK} \phi \subseteq \mathsf{CK} (\neg \mathsf{CK} \phi)$ .

*Proof.* Assume  $(\neg \mathsf{CK} \phi) w$  and, with Theorem 4.4.2,  $w \propto_{\mathsf{Agent}} v$ . To prove the negation  $(\neg \mathsf{CK} \phi) v$ , which means  $\mathsf{CK} \phi v \to \mathbb{O}$ , we may assume  $\mathsf{CK} \phi v$  to reach a contradiction. By symmetry, we have  $v \propto_{\mathsf{Agent}} w$ , then Lemma 4.4.1 allows us to derive  $\mathsf{CK} \phi w$ , which contradicts the first assumption.

Theorem 4.4.7. CK is a knowledge operator.

*Proof.* The full theorem is formed by combining Lemmas 4.4.3 to 4.4.6.  $\Box$ 

# Chapter V Coalition Logic

In this chapter, we introduce coalition logic and its semantics based on game forms and playable effectivity functions. As in previous chapters, this will at first be in isolation from type theory, using traditional set-theoretic notions to define the concepts of the logic. We will go on to describe a formulation of these concepts appropriate for use in type theory, and sketch an equivalence proof for game forms and playable effectivity functions.

The content of this chapter is based on work in progress, without a corresponding Agda formalisation in Appendix A. It would take significant future effort to develop the following outline so it can be fully formalised in a proof assistant.

#### 5.1 Background

Coalition logic is a modal logic introduced by Marc Pauly [49]. Like epistemic logic, it is an agent-based logic, the fundamental concept being the *coalition*, a subset of agents. Rather than considering each agents in isolation, coalition logic explores what agents can achieve when they act in unison.

In the syntax of the logic, each coalition C has a modal operator,  $[C] \phi$ , with the reading that the agents of C, if they so are inclined to work together, could force  $\phi$  to hold after their collective action, regardless of any other agents' actions. The base logic has no concept of knowledge, essentially modelling perfectinformation strategic games, but it can be extended with modal operators for knowledge to model epistemic scenarios [1, 56].

The semantics of coalition logic is given by game forms. These are descriptions of a step in some strategic game where players may act simultaneously. A game form consists of a finite, non-empty set of agents, N, a set of actions,  $\mathcal{A}_i$  for each  $i \in N$ , a set of outcomes, O, and a function which maps a choice of action for each agent to an outcome,  $o: (\prod_{i \in N} \mathcal{A}_i) \to O$ . We can derive an *effectivity function* for each game form G. An effectivity function  $E_G : \mathcal{P}(N) \to \mathcal{P}(\mathcal{P}(O))$  describes, for each coalition of agents, those subsets of outcomes for which they are effective. That is, if  $X \in E_G(C)$ , coalition C has a joint strategy which guarantees the actual outcome must come from X, regardless of the strategies chosen by the other agents not in C. To define the function precisely, some additional notation is introduced.

Let *C* range over coalitions, subsets of agents,  $C \in \mathcal{P}(N)$ , and *X* range over subsets of outcomes,  $X \in \mathcal{P}(O)$ . By  $\overline{C}$  and  $\overline{X}$ , we refer to the complements of *C* and *X* relative to their full sets,  $\overline{C} = N \setminus C$ ,  $\overline{X} = O \setminus X$ . A strategy profile,  $\sigma_C$ , is a tuple containing the choice of action for each agent in a coalition,  $\sigma_C : \prod_{(i \in C)} \mathcal{A}_i$ . A global strategy profile,  $\sigma : \prod_{(i \in N)} \mathcal{A}_i$ , can be constructed by combining strategy profiles of complementary coalitions,  $\sigma = \sigma_C \oplus \sigma_{\overline{C}}$ . The effectivity function for game form *G* may then be defined.

$$E_G(C) = \{ X \mid \exists_{\sigma_C} \forall_{\sigma_{\overline{C}}} \ o(\sigma_C \oplus \sigma_{\overline{C}}) \in X \}$$

Is is convenient to work abstractly with effectivity functions rather than game form definitions. Pauly shows that there is a set of properties for effectivity functions, *playability*, which characterises when an arbitrary effectivity function E is the effectivity function  $E_G$  for some game form G.

- SAFETY,  $\emptyset \notin E(C)$
- LIVENESS,  $O \in E(C)$
- *N*-Maximality,  $\overline{X} \notin E(\emptyset) \Rightarrow X \in E(N)$
- OUTCOME MONOTONICITY,  $X_1 \subseteq X_2 \Rightarrow X_1 \in E(C) \Rightarrow X_2 \in E(C)$
- SUPERADDITIVITY,  $C_1 \cap C_2 = \emptyset \Rightarrow X_1 \in E(C_1) \Rightarrow X_2 \in E(C_2) \Rightarrow X_1 \cap X_2 \in E(C_1 \cup C_2)$

Two further properties can be shown to follow from those listed above.

- REGULARITY,  $X \in E(C) \Rightarrow \overline{X} \notin E(\overline{C})$
- Coalition Monotonicity,  $C_1 \subseteq C_2 \Rightarrow X \in E(C_1) \Rightarrow X \in E(C_2)$

Given a game form G, the playability of  $E_G$  is a routine question of checking each of the properties. The inverse direction, starting from a playable effectivity function E, requires construction of a game form G such that E and  $E_G$  are equivalent. For the rest of this chapter, we work towards giving a type-theoretic formulation for this direction of the correspondence.

#### 5.2 Game Form Construction

Our construction of the game form is based on the idea of Pauly's original construction [49]. However, this is heavily set-theoretic and poses significant technical challenges to formulation in type theory. This section will first give an outline of the construction, and we will discuss some of the technical details in later sections.

Given a playable effectivity function  $E : \mathcal{P}(N) \to \mathcal{P}(\mathcal{P}(O))$ . The constructed game form will use the same agent and outcome sets, so we just need to define the actions available to each agent and the outcome function.

An action for agent *i* consists of a choice of coalition *C* such that  $i \in C$ , a set of outcome states  $X \in E(C)$ , a specific outcome  $x \in X$ , and a natural number. Formally,  $\mathcal{A}_i$  is defined as the following set.

$$\mathcal{A}_i = \{ \langle C, X, x, t \rangle \mid C \in \mathcal{P}(N), i \in C, X \in E(C), x \in X, t \in \mathbb{N} \}$$

Intuitively, if an agent chooses action  $\langle C, X, x, t \rangle$ , it means they want to be part of coalition C, cooperating to achieve an outcome in X, and they have a personal preference for specific outcome x if the choice were up to them. The natural number t will be used in the outcome function to determine which agent gets to make the final decision.

Given a strategy profile  $\sigma : \Pi_{(i \in N)} \mathcal{A}_i$ , we have a choice of action for each agent,  $\langle C_i, X_i, x_i, t_i \rangle$ . We call C a  $\sigma$ -cooperative coalition when, for every agent  $i \in C$ ,  $C = C_i$ , and for every pair of agents  $i, j \in C$ ,  $X_i = X_j$ . That is, a  $\sigma$ -cooperative coalition C agrees to work together towards a common set of outcomes  $X_C$ .

Let  $[C_1, \ldots, C_m]$  be the list of all the non-empty  $\sigma$ -cooperative coalitions. There must be finitely many since N is itself finite. Let  $C_0$  be the coalition of agents not in a  $\sigma$ -cooperative coalition, partitioning the full set of agents  $[C_0, C_1, \ldots, C_m]$ . Since  $C_0$  has no agreed set of outcomes to aim for, define  $X_{C_0} = O$ , and define  $X_{\sigma}$  as the intersection of all chosen outcome sets  $X_{C_k}$ .

$$X_{\sigma} = \bigcap_{k=0}^{m} X_{C_k} = \bigcap_{k=1}^{m} X_{C_k}$$

It cannot be the case that  $X_{\sigma} = \emptyset$ . Since  $[C_0, C_1, \ldots, C_m]$  is a partition, and all  $X_{C_k} \in E(C_k)$  by definition of  $\mathcal{A}$ , the superadditivity property of E implies that  $X_{\sigma} \in E(N)$ , and safety tells us  $\emptyset \notin E(N)$ .

The outcome function will be defined to choose an outcome  $x \in X_{\sigma}$ . Recall that the number of agents is finite, and let n be their number, n = |N|. Labelling the set of agents by natural numbers  $\{0, \ldots, n-1\}$ , we sum modulo n the choice

of  $t_i$  for each agent,  $d = t_0 + \cdots + t_{n-1} \mod n$ , to define the agent who decides the outcome.

It must be the case that  $x_d \in X_{C_d}$ , but it is not guaranteed that  $x_d \in X_{\sigma}$ . If  $x_d \notin X_{\sigma}$ , we revert to an arbitrary choice function  $H : \prod_{(X \in E(N))} X$ , which exists constructively due to the safety property,  $\emptyset \notin E(N)$ . Since  $X_{\sigma} \in E(N)$  from the above superadditivity argument, we can define the outcome function.

$$o(\sigma) = \begin{cases} x_d & \text{if } x_d \in X_\sigma \\ H(X_\sigma) & \text{otherwise} \end{cases}$$

Thus we have constructed the game form G. To show the correspondence between game forms and playable effectivity functions, we need to demonstrate that E is equivalent to the derived effectivity function  $E_G$ .

Lemma 5.2.1. If  $X \in E(C)$ , then  $X \in E_G(C)$ .

Proof. Assume  $X \in E(C)$ , then we have to show  $\exists_{\sigma_C} \forall_{\sigma_{\overline{C}}} o(\sigma_C \oplus \sigma_{\overline{C}}) \in X$ . We define  $\sigma_C$  by giving an action  $\langle C_i, X_i, x_i, t_i \rangle$  for every  $i \in C$ . Set  $C_i = C$  and  $X_i = X$ , choosing  $x_i$  and  $t_i$  arbitrarily — our assumption that  $X \in E(C)$  together with the safety property,  $\emptyset \notin E(C)$ , ensure an  $x \in X$  can be chosen.

Given any counter strategy  $\sigma_{\overline{C}}$ , we have strategy profile  $\sigma = \sigma_C \oplus \sigma_{\overline{C}}$ , where C is  $\sigma$ -cooperative by definition. C is therefore one of the partition classes whose chosen outcome set X is used to define  $X_{\sigma}$  by intersection, so  $X_{\sigma} \subseteq X$ . Since  $o(\sigma) \in X_{\sigma}$ , this subset relation implies  $o(\sigma) \in X$ , as desired.

Lemma 5.2.2. If  $X \in E_G(C)$ , then  $X \in E(C)$ .

*Proof.* We assume that coalition membership is decidable, so we distinguish two cases, when C = N and when there is at least one  $i \in N$  where  $i \in \overline{C}$ .

In the first case, we assume  $X \in E_G(N)$ . By regularity of  $E_G$ , this implies  $\overline{X} \notin E_G(\emptyset)$ . Using the contrapositive of Lemma 5.2.1, we can derive  $X \notin E(\emptyset)$ . Then, N-maximality of E reaches the desired conclusion  $X \in E(N)$ .

Now consider the case with at least one  $i \in \overline{C}$ . Assume  $X \in E_G(C)$ , whose definition means that C has a strategy profile  $\sigma_C$  such that  $\forall_{\sigma_{\overline{C}}} o(\sigma_C \oplus \sigma_{\overline{C}}) \in X$ . Instantiating this with a strategy profile  $\sigma_{\overline{C}}$  where  $C_i = N$  and  $X_i = O$  for every  $i \in \overline{C}$  makes  $X_{\sigma}$  depend only on  $\sigma_C$ .

$$X_{\sigma} = \bigcap \{ X_{C'} \mid C' \subseteq C, \ C' \text{ is } \sigma \text{-cooperative} \}$$

Since, by definition of  $\mathcal{A}$ , all  $X_{C'} \in E(C')$ , superadditivity and coalition monotonicity imply  $X_{\sigma} \in E(C)$ . If we can give appropriate values for  $x_i$  and  $t_i$ for each  $i \in \overline{C}$  so that  $X_{\sigma} \subseteq X$ , this will show the desired  $X \in E(C)$  by outcome monotonicity. For each  $x \in X_{\sigma}$ , since we have at least one  $i \in \overline{C}$ , we can set  $x_i = x$ , and tweak the values of  $t_i$  so that  $d = t_0 + \cdots + t_{n-1} \mod n$  ends up as one of i. This ensures  $o(\sigma_C \oplus \sigma_{\overline{C}}) = x$ , and so  $x \in X$  by our earlier assumption, therefore  $X_{\sigma} \subseteq X$ .

**Theorem 5.2.3.** E and  $E_G$  are equivalent,  $E(C) \subseteq E_G(C)$  and  $E(C) \supseteq E_G(C)$ .

*Proof.* The left-to-right direction is Lemma 5.2.1, while the right-to-left direction is Lemma 5.2.2.  $\hfill \Box$ 

### 5.3 Decidable Subsets

Back in the setting of type theory, we seek to represent the concepts of coalition logic. We begin by assuming appropriate types for agents and outcomes.

Unlike in Chapter IV, we impose the restriction that the agent type is finite, as well as non-empty. We fix a non-zero natural number n as the number of agents,  $\Sigma_{(n:\mathbb{N})} n \not\equiv \text{zero}$ , and define Agent = Fin n: Type. Our agents therefore correspond to the natural numbers  $\{0, \ldots, n-1\}$ , and since n must be a successor, we can always construct the agent  $\text{zero}_{\text{Fin}}$ : Agent. We assume Outcome : Type without any restrictions.

We must develop a library to work with subsets of these types. A type family  $T \rightarrow \mathsf{Type}$ , representing predicates over type T, can be seen as a subset of the type's elements — a term t : T is part of the subset if and only if type T t is inhabited. However, in the constructive setting of type theory, there is no guarantee that predicates are decidable in general, which will make working with this representation of subsets difficult to impossible. We instead encode the type of *decidable subsets* by functions into the Booleans,  $T \rightarrow 2$ , which correspond to decidable predicates whose truth values always reduce to true or false.

**Definition 5.3.1.** A decidable subset of type T is a function mapping from T into the Booleans.

$${\mathcal P}: \mathsf{Type} o \mathsf{Type}$$
  
 ${\mathcal P}: T = T o 2$ 

Given a decidable subset  $S : \mathcal{P} T$  and a term t : T, we obtain the truth value of subset membership by function application, S t : 2. From here, the notation of set theory can be defined using Boolean algebra, with the true and false constants, and the basic operators  $\mathsf{not}_2$ ,  $\mathsf{or}_2$ , and  $\mathsf{and}_2$ .

Definition 5.3.2. Set-theoretic notation defined for decidable subsets.

$$\begin{aligned} \mathsf{full} &: \Pi_{(T:\mathsf{Type})} \ \mathcal{P} \ T & & \varnothing : \Pi_{(T:\mathsf{Type})} \ \mathcal{P} \ T \\ \mathsf{full} &= \lambda t. \ \mathsf{true} & & & \varnothing &= \lambda t. \ \mathsf{false} \end{aligned}$$

$\overline{S} : \Pi_{(T:Type)} \mathcal{P} T \to \mathcal{P} T$ $\overline{S} = \lambda t. not_2 (S t)$	$ \langle : \Pi_{(T:Type)} \mathcal{P} T \to \mathcal{P} T \to \mathcal{P} T \\ A \setminus B = \lambda t. \text{ and}_2 (A t) (not_2 (B t)) $
$\cup: \Pi_{(T:Type)} \mathcal{P} \ T \to \mathcal{P} \ T \to \mathcal{P} \ T$ $A \cup B = \lambda t. \text{ or}_2 \ (A \ t) \ (B \ t)$	$T \qquad \cap: \Pi_{(T:Type)} \mathcal{P} \ T \to \mathcal{P} \ T \to \mathcal{P} \ T$ $A \cap B = \lambda t. \text{ and}_2 \ (A \ t) \ (B \ t)$
$\in:\Pi_{(T:Type)} \ T \to \mathcal{P} \ T \to Type$	$\subseteq : \Pi_{(T:Type)} \ \mathcal{P} \ T \to \mathcal{P} \ T \to Type$
$t\in S=S \ t\equiv {\rm true}$	$A \subseteq B = \Pi_{(t:T)} \ t \in A \to t \in B$

Most of these definitions are standard, needing little explanation. To state the set membership  $t \in S$  as a proposition in type theory, that is as a type, we assert that the truth value resulting from S t is propositionally equal to true. The statement  $t \notin S$  can be defined simply as  $t \in S \to \mathbb{O}$ .

**Theorem 5.3.3.** Set membership is decidable,  $t \in S + t \notin S$ .

*Proof.* The full statement is  $S t \equiv \text{true} + (S t \equiv \text{true} \rightarrow 0)$ . By case analysis on the Boolean value S t, using the eliminators  $\text{elim}_+$  and  $\text{elim}_1$ , we have to demonstrate the statement when S t is true and when it is false.

In the true case, we choose inl, proving the left equality,  $refl_{\equiv}$  true. In the false case, we choose inr, showing that false  $\equiv$  true leads to absurdity, which is implied by Theorem 2.5.9.

By pattern matching on the result of Theorem 5.3.3, it is possible to derive a function which decides the statement  $t \in S$ , giving the result as a Boolean value, inl  $l \mapsto \text{true}$ , inr  $r \mapsto \text{false}$ . This result of this function will coincide with the value of S t.

We do not define equality of decidable subsets with a dedicated relation. It could be stated  $(A \subseteq B) \times (B \subseteq A)$ , as is usual in set theory. However, as decidable subsets are functions, it is more convenient to postulate the principle of function extensionality, funExt :  $(\Pi_{(t:T)} A t \equiv B t) \rightarrow A \equiv B$ . We consider two decidable subsets propositionally equal when they are pointwise equal, giving the same truth values for set membership at all indices.

For example, to show that  $\overline{\mathsf{full}} \equiv \emptyset$ , by appealing to function extensionality we can instead show  $\Pi_{(t:T)}$   $\overline{\mathsf{full}}$   $t \equiv \emptyset$  t. This reduces the question of subset equality into a Boolean equality,  $\mathsf{not}_2$  true  $\equiv$  false, which can be shown with the usual proof methods of type theory,  $\mathsf{funExt}$  ( $\lambda t$ .  $\mathsf{refl}_{\equiv}$  false).

**Lemma 5.3.4.** Complementary subsets are disjoint,  $S \cap \overline{S} \equiv \emptyset$ .

*Proof.* funExt ( $\lambda t$ . elim<sub>+</sub> (elim<sub>1</sub> (refl<sub>=</sub> false)) (elim<sub>1</sub> (refl<sub>=</sub> false)) (S t))

Using function extensionality, this can be proved as an instance of the Boolean theorem  $\Pi_{(b:2)}$  and b (not b)  $\equiv$  false, where b is instantiated with the set membership truth value S t. The proof of the Boolean theorem is given by two cases, when b is true and when b is false, both of which reduce to proving false  $\equiv$  false.  $\Box$ 

**Lemma 5.3.5.** Dual to the above,  $S \cup \overline{S} \equiv \mathsf{full}$ .

*Proof.* funExt ( $\lambda t$ . elim<sub>+</sub> (elim<sub>1</sub> (refl<sub>=</sub> true)) (elim<sub>1</sub> (refl<sub>=</sub> true)) (S t))

The only difference from the proof of Lemma 5.3.4 is that the two cases reduce to true  $\equiv$  true rather than false  $\equiv$  false.

Using propositional equality for subsets gives us access to the full suite of tools for equality proofs, including the powerful congruence substitution principle. We will need many more properties of set equality, which can be proved first as Boolean equalities and then elevated to set equalities by function extensionality.

Given this representation of subsets, we are able to state the type of effectivity functions for our type of agents and outcomes,  $\mathcal{P} \operatorname{Agent} \to \mathcal{P} (\mathcal{P} \operatorname{Outcome})$ . We can further define the type of playable effectivity functions.

**Definition 5.3.6.** A playable effectivity function is an effectivity function satisfying the playability properties.

$$\begin{array}{l} \mathsf{Playable}:\mathsf{Type}\\ \mathsf{Playable}=&\Sigma_{(E:\mathcal{P}\,\mathsf{Agent}\to\mathcal{P}\,(\mathcal{P}\,\,\mathsf{Outcome}))}\\ &\Pi_{(C:\mathcal{P}\,\mathsf{Agent})} \ \varnothing \not\in E \ C\\ &\times\Pi_{(C:\mathcal{P}\,\mathsf{Agent})} \ \mathsf{full}\in E \ C\\ &\times\Pi_{(X:\mathcal{P}\,\mathsf{Outcome})} \ \overline{X} \not\in E \ \varnothing \to X \in E \ \mathsf{full}\\ &\times\Pi_{(C:\mathcal{P}\,\mathsf{Agent})} \ \Pi_{(X_1,X_2:\mathcal{P}\,\mathsf{Outcome})} \ X_1 \subseteq X_2 \to X_1 \in E \ C \to X_2 \in E \ C\\ &\times\Pi_{(C_1,C_2:\mathcal{P}\,\mathsf{Agent})} \ \Pi_{(X_1,X_2:\mathcal{P}\,\mathsf{Outcome})} \ C_1 \cap C_2 \equiv \varnothing \to\\ &X_1 \in E \ C_1 \to X_2 \in E \ C_2 \to X_1 \cap X_2 \in E \ (C_1 \cup C_2) \end{array}$$

The additional properties of regularity and coalition monotonicity can be shown to follow from playability.

**Theorem 5.3.7.** Playable effectivity functions are regular,  $X \in E \ C \to \overline{X} \notin E \ \overline{C}$ .

*Proof.* Assume  $X \in E \ C$  and, towards a contradiction,  $\overline{X} \in E \ \overline{C}$ . Since C and  $\overline{C}$  can be shown to be disjoint, superadditivity can be used to conclude  $X \cap \overline{X} \in E \ (C \cup \overline{C})$ . By using congruence to substitute with proofs of the equalities  $X \cap \overline{X} \equiv \emptyset$  and  $C \cup \overline{C} \equiv$  full, we have  $\emptyset \in E$  full, contradicting the safety property of playability.

**Theorem 5.3.8.** Playable effectivity functions are coalition monotonic,  $C_1 \subseteq C_2 \rightarrow E \ C_1 \subseteq E \ C_2$ .

*Proof.* Assume  $C_1 \subseteq C_2$  and  $X \in E$   $C_1$ , with the goal of proving  $X \in E$   $C_2$ . Define  $C_3 = C_2 \setminus C_1$ , which gives equalities  $C_1 \cup C_3 \equiv C_2$  and  $C_1 \cap C_3 \equiv \emptyset$ . By the liveness property of playability, full  $\in E \ C_3$ , so superadditivity implies  $X \cap \mathsf{full} \in E \ (C_1 \cup C_3)$ . Substituting, using congruence with proofs of  $X \cap \mathsf{full} \equiv X$  and  $C_1 \cup C_3 \equiv C_2$ , gives the desired conclusion  $X \in E \ C_2$ .

These proofs show just some of the standard set-theoretic properties that are needed. As described earlier, these follow from function extensionality and basic Boolean properties. Most of these properties will not be proved explicitly, but we will give one final example which was used in Theorem 5.3.8.

#### **Lemma 5.3.9.** $A \subseteq B \rightarrow A \cup (B \setminus A) \equiv B$

*Proof.* Assume  $A \subseteq B$ . By function extensionality, we can reduce the conclusion to the Boolean statement  $\operatorname{or}_2(A \ t) (\operatorname{and}_2(B \ t) (\operatorname{not}_2(A \ t))) \equiv B \ t$ . By case analysis, this can be checked to hold except when  $A \ t \equiv \operatorname{true}$  and  $B \ t \equiv \operatorname{false}$ .

However, instantiating the assumption of  $A \subseteq B$  at element t, it is the Boolean statement  $A \ t \equiv \mathsf{true} \to B \ t \equiv \mathsf{true}$ . Therefore, by congruence, we can construct  $\mathsf{true} \equiv \mathsf{false}$  in the problematic case, ruling it absurd by Theorem 2.5.9, completing the proof using  $\mathsf{elim}_0$ .

## 5.4 Partial Functions

We can define a type for game forms using our existing Agent and Outcome types.

**Definition 5.4.1.** A game form is a dependent pair  $\langle \mathcal{A}, o \rangle$  where  $\mathcal{A}$  is a type family indexed by agents, and o is a function assigning an outcome to a choice of action for each agent.

$$\begin{split} &\mathsf{GameForm}:\mathsf{Type}\\ &\mathsf{GameForm}=\Sigma_{(\mathcal{A}:\mathsf{Agent}\to\mathsf{Type})}\ (\Pi_{(i:\mathsf{Agent})}\ \mathcal{A}_i)\to\mathsf{Outcome} \end{split}$$

The outcome function o is total, only assigning outcomes to strategy profiles for the full set of agents,  $\sigma : \Pi_{(i:\text{Agent})} \mathcal{A}_i$ . However, the definition of effectivity function  $E_G$  relies on having strategy profiles for coalitions of agents,  $\sigma_C$ ,  $\sigma_{\overline{C}}$ . We use partial functions to bridge this gap.

**Definition 5.4.2.** A partial function mapping from A: Type to an element of type family  $B: A \to \text{Type}$  is a dependent function  $\prod_{(a:A)} B a + 1$ .

Of course, strictly speaking,  $A \rightarrow B$  is not a partial function, but a total function into a codomain with exactly one additional element. The intuition is

that if a partial function maps  $a \mapsto b$ , this is represented by mapping  $a \mapsto \text{inl } b$ . When the partial function does not provide a mapping from a, we represent this by mapping to the dummy element,  $a \mapsto \text{inr } \star$ .

It can also be the case by this definition that a so-called partial function is actually total, mapping to  $a \mapsto \text{inl} b$  for all a : A. We can characterise elements of A for which there are mappings, the domain of the partial function, using our decidable subset type.

**Definition 5.4.3.** The domain of a partial function  $A \rightarrow B$  is the subset of A for which the function has mappings.

Domain : 
$$\Pi_{(A:\mathsf{Type})} \Pi_{(B:A \to \mathsf{Type})} (A \rightharpoonup B) \to \mathcal{P} A$$
  
Domain  $f = \lambda a.$  elim<sub>+</sub> ( $\lambda b.$  true) (elim<sub>1</sub> false) ( $f a$ )

For example, the partial function with no mappings has an empty domain, Domain  $(\lambda(a : A). \text{ inr } \star) \equiv \emptyset$ . Given a total function  $f : \prod_{(a:A)} B a$ , we can construct a corresponding partial function, Domain  $(\lambda(a : A). \text{ inl } (f a)) \equiv \text{full}.$ 

We can join the domains of two partial functions  $f, g : A \rightarrow B$  produce another partial function  $f \oplus g : A \rightarrow B$ . The resulting function only maps  $a \mapsto \text{inr} \star$  when neither joined function has a mapping from a.

$$(f \oplus g) \ a = \begin{cases} \text{inl } b \ \text{if } f \ a \equiv \text{inl } b \\ \text{inl } b \ \text{if } g \ a \equiv \text{inl } b \\ \text{inr } \star \text{ otherwise} \end{cases}$$

While  $\oplus$  is always associative, it is only commutative under the assumption that the domains of f and g are disjoint. This is because we default to the mapping provided by f in the case where both f and g provide mappings. The join produces a function whose domain is the union of the two original functions, Domain  $(f \oplus g) \equiv$  Domain  $f \cup$  Domain g.

Given a game form  $G = \langle \mathcal{A}, o \rangle$ , a strategy profile for a coalition C is a partial function  $\sigma_C : \text{Agent} \rightarrow \mathcal{A}$ , mapping to an action in  $\mathcal{A}_i$  if and only if  $i \in C$ . That is, Domain  $\sigma_C \equiv C$ . Since Lemma 5.3.5 tells us that  $C \cup \overline{C} \equiv \text{full}$ , given strategy profiles  $\sigma_C$  and  $\sigma_{\overline{C}}$  with these domains, we can show Domain  $(\sigma_C \oplus \sigma_{\overline{C}}) \equiv \text{full}$ . From this, we are able to construct a total strategy profile  $\sigma : \Pi_{(i:\text{Agent})} \mathcal{A}_i$ .

**Lemma 5.4.4.** Given a partial function  $f : A \rightarrow B$  and an a : A such that  $a \in \text{Domain } f$ , then we can construct an element of type B a.

*Proof.* Assume that  $a \in \text{Domain } f$ , which means Domain  $f a \equiv \text{true.}$  Apply f a, which gives an element of type B a + 1. Pattern matching on the result, this must be of the form inl b where b : B a, or inr  $\star$ .

In the first case, we can simply return b. In the second case, we cannot construct an element of type B a, but Domain f a must have been false, refining our assumption to false  $\equiv$  true. By Theorem 2.5.9 and elim<sub>0</sub>, this case can be dismissed as absurd.

**Theorem 5.4.5.** When the domain of a partial function is the full set, a total function can be constructed,  $\Pi_{(f:A \rightarrow B)}$  Domain  $f \equiv \text{full} \rightarrow \Pi_{(a:A)} B a$ .

*Proof.* Assume that Domain  $f \equiv$  full. Given an a : A, it must be the case that  $a \in$  full by definition. By congruence using our assumed equality, it must also be the case that  $a \in$  Domain f. Our element of type B a can then be constructed through Lemma 5.4.4.

#### 5.5 Future Work

We have given suitable constructions in type theory to represent game forms and playable effectivity functions. However, a full formalisation of the construction outlined in Section 5.2 is incomplete, and is left as future work. This section discusses some of the remaining challenges of the formalisation.

First, when defining the transformation from playable effectivity function E to game form, the definition  $\mathcal{A}$  is clear, as a type family indexed by Agent.

**Definition 5.5.1.** An action for agent i is a choice of coalition that i is part of, a set of outcomes that coalition is effective for, a specific outcome from that set, and a natural number.

$$\begin{aligned} \mathcal{A} : \mathsf{Agent} &\to \mathsf{Type} \\ \mathcal{A}_i = &\Sigma_{(C:\mathcal{P} \, \mathsf{Agent})} & i \in C \\ &\times &\Sigma_{(X:\mathcal{P} \, \mathsf{Outcome})} & X \in E \ C \\ &\times &\Sigma_{(x:\mathsf{Outcome})} & x \in X \\ &\times &\mathbb{N} \end{aligned}$$

With projection functions  $\mathcal{A}_C$  and  $\mathcal{A}_X$ , we can then define the  $\sigma$ -cooperative predicate for coalitions relative to a specific strategy profile  $\sigma$ .

**Definition 5.5.2.** Given a strategy profile  $\sigma$ , a coalition C is  $\sigma$ -cooperative when all of its members choose actions in  $\sigma$  which agree on C and X.

$$\begin{aligned} \operatorname{coop} &: (\Pi_{(i:\operatorname{Agent})} \ \mathcal{A}_i) \to \mathcal{P} \ \operatorname{Agent} \to \operatorname{Type} \\ \operatorname{coop} \ \sigma \ C &= \Pi_{(i:\operatorname{Agent})} \ i \in C \to C \equiv \mathcal{A}_C \ (\sigma \ i) \\ &\times \Pi_{(j:\operatorname{Agent})} \ j \in C \to \mathcal{A}_X \ (\sigma \ i) \equiv \mathcal{A}_X \ (\sigma \ j) \end{aligned}$$

The partitioning process will require a proof that this predicate is decidable, coop  $\sigma C + (\text{coop } \sigma C \to \mathbb{O})$ . One possibility for a data structure to describe the partition would be an inductive list.

Starting from a list of all agents, we can map a function pointwise which pairs the  $\langle C_i, X_i \rangle$  choices for each agent, then filter out those elements whose  $C_i$  is not  $\sigma$ -cooperative. This gives a list of all non-empty,  $\sigma$ -cooperative coalitions and their choice of  $X_C$ , and we must prove each pair of coalitions disjoint so we may later use superadditivity. To define  $X_{\sigma}$ , we can use a recursive function which intersects the  $X_i$  choices of the partition from a base case of full.

 $\begin{array}{ll} \text{intersect}: \mathsf{List} \ (\mathcal{P} \ \mathsf{Agent} \times \mathcal{P} \ \mathsf{Outcome}) \to \mathcal{P} \ \mathsf{Outcome} \\ \text{intersect} \ [] &= \mathsf{full} \\ \text{intersect} \ (\langle C_i, \ X_i \rangle :: t) = X_i \cap \text{intersect} \ t \end{array}$ 

However, this must be done in such a way that the proof terms which verify that  $X_i \in E \ C_i$  for each *i*, which are included the structure of action  $\sigma \ i : \mathcal{A}_i$ , are retained throughout the process. At each recursive step, we need to use superadditivity with this proof to produce a proof of  $X_{\sigma} \in E$  full along with  $X_{\sigma}$ . This will justify  $X_{\sigma} \neq \emptyset$  by safety, and allow for the definition of outcome function *o*.

With a game form  $G = \langle \mathcal{A}, o \rangle$ , we have shown in Theorem 5.4.5 that two complementary strategy profiles can be joined to give a complete strategy profile  $\sigma : \Pi_{(i:\text{Agent})} \mathcal{A}_i$ . Using this, we can formulate a version of the following settheoretic statement in type theory, to be used in the definition of effectivity function  $E_G$ .

$$\exists \sigma_C \forall \sigma_{\overline{C}} \ o(\sigma_C \oplus \sigma_{\overline{C}}) \in X$$

To actually define the function, GameForm  $\rightarrow \mathcal{P}$  Agent  $\rightarrow \mathcal{P}$  ( $\mathcal{P}$  Outcome), the result is a decidable subset. When applied to a coalition  $C : \mathcal{P}$  Agent and a set of outcomes  $X : \mathcal{P}$  Outcome, this means the result must be Boolean,  $E_G C X : 2$ . However, the obvious translation into type theory, using a  $\Sigma$ -type and a  $\Pi$ -type to express the quantifiers, and omitting details establishing the domains of the partial functions, is a statement residing in a type universe.

$$\Sigma_{(\sigma_C: \operatorname{Agent} \to \mathcal{A})} \prod_{(\sigma_{\overline{C}}: \operatorname{Agent} \to \mathcal{A})} o(\sigma_C \oplus \sigma_{\overline{C}}) \in X : \operatorname{Type}$$

To define the effectivity function with the type given above, it must be shown that this statement is decidable. Alternatively, the type of  $E_G$  could be reformulated as a non-decidable predicate,  $\mathcal{P}$  Agent  $\rightarrow \mathcal{P}$  Outcome  $\rightarrow$  Type, but then the notion of playability in Definition 5.3.6 must also be reformulated.

# Chapter VI Concluding Remarks

In Chapter III of this thesis, we developed an embedding for epistemic modal logic within type theory. At first, this was without a type of agents, and we defined what it means to be a knowledge operator in this setting, giving a set of properties such an operator must satisfy. This included the preservation of semantic entailment property, which is an infinitary deduction rule that subsumes two axioms of the finitary axiomatic system of epistemic logic. We proved this knowledge operator semantics, including the additional rule, coincides with the traditional semantics based on equivalence relations, establishing an isomorphism between the two.

In Chapter IV, we extended the embedding with a type of agents, giving each a knowledge relation, and therefore a knowledge operator by using the isomorphism of Chapter III. We introduced a universal knowledge operator to the embedding, and used it to define a common knowledge operator as a coinductive predicate, whose proofs are infinite data structures. We proved that coinductive common knowledge is equivalent to a type family which iterates universal knowledge, and also to the relational interpretation of common knowledge, by taking the union and transitive closure of the individual agents' relations. Formalisations of all results from Chapters III and IV are available in Appendix A, verified by the Agda proof assistant.

In Chapter V, we turned our attention to coalition logic, and its semantics based on game forms and playable effectivity functions. We outlined an equivalence proof for these with the intention of formalisation in type theory. Type-theoretic representations are given, using a type of decidable subsets to define playable effectivity functions, and partial functions used to solve a technical problem when defining an effectivity function from a given game form. However, the formalisation is not currently in a state that can be verified by proof assistant, and completion of the formalisation is left as future work. We sketched the remaining challenges which need to be overcome for the full formalisation.

# Appendix A

# **Agda Formalisations**

Formalisations based on the content of Chapter II, Chapter III, and Chapter IV, developed in the Agda proof assistant [48].

# A.1 Dependent Type Theory

```
{-
   AGENT-BASED LOGICS IN DEPENDENT TYPE THEORY
   - by Colm Baston
   CHAPTER II: DEPENDENT TYPE THEORY
   - checked with Agda version 2.7.0.1
-}
```

module TypeTheory where

-- levels for universe polymorphism open import Agda.Primitive variable  $\alpha \ \beta \ \gamma$ : Level

{- BASIC TYPE DEFINITIONS -}

```
data 0 : Set where
  -- no constructors
```

```
\mathsf{elim} \cdot \mathbb{O} : \{ T : \mathbb{O} \to \mathsf{Set} \; \alpha \} \; (e : \mathbb{O}) \to \; T \; e\mathsf{elim} \cdot \mathbb{O} \; ()
```

data 1 : Set where \*:1  $\mathsf{elim}\text{-}1: \{T: 1 \to \mathsf{Set} \; \alpha\} \to T \star \to (u: 1) \to T \; u$ elim-1  $t \star = t$ data \_+\_ (A : Set  $\alpha$ ) (B : Set  $\beta$ ) : Set ( $\alpha \sqcup \beta$ ) where inl :  $A \rightarrow A + B$ inr :  $B \rightarrow A + B$ elim+ : {A : Set  $\alpha$ } {B : Set  $\beta$ } {T :  $A + B \rightarrow$ Set  $\gamma$ }  $\rightarrow$  $((a: A) \rightarrow T (\mathsf{inl} a)) \rightarrow$  $((b:B) \rightarrow T (\text{inr } b)) \rightarrow$  $(s: A + B) \rightarrow T s$ elim-+ f (inl a) = f aelim+ g (inr b) = g b2 : Set 2 = 1 + 1true: 2 true = inl  $\star$ false : 2  $false = inr \star$ data  $\Sigma$  (A : Set  $\alpha$ ) (B :  $A \rightarrow Set \beta$ ) : Set ( $\alpha \sqcup \beta$ ) where  $\_,\_: (a: A) \to B \ a \to \Sigma \ A \ B$  $\_\times\_$ : Set  $\alpha \to$ Set  $\beta \to$ Set  $(\alpha \sqcup \beta)$  $A \times B = \Sigma A (\lambda \_ \rightarrow B)$  $\mathsf{elim}\text{-}\Sigma: \{A: \mathsf{Set}\ \alpha\}\ \{B: A \to \mathsf{Set}\ \beta\}\ \{T: \Sigma\ A\ B \to \mathsf{Set}\ \gamma\} \to$  $((a:A) (b:B a) \rightarrow T (a,b)) \rightarrow$  $(p: \Sigma A B) \rightarrow T p$ elim- $\Sigma f(a, b) = f a b$ 

record EquivRel (A : Set  $\alpha$ ) : Set ( $\alpha \sqcup$  lsuc  $\beta$ ) where field  $\mathsf{R} : A \to A \to \mathsf{Set} \ \beta$ refl :  $\{a : A\} \rightarrow a \mathbb{R} a$ sym :  $\{a \ b : A\} \rightarrow a \ \mathsf{R} \ b \rightarrow b \ \mathsf{R} \ a$ trans :  $\{a \ b \ c : A\} \rightarrow a \ \mathsf{R} \ b \rightarrow b \ \mathsf{R} \ c \rightarrow a \ \mathsf{R} \ c$ {- PROPOSITIONAL LOGIC EXAMPLES -} module Propositional ( $P \ Q \ R : Set$ ) where  $\mathsf{refl} \to : P \to P$ refl- $\rightarrow p = p$ trans- $\rightarrow$ :  $(P \rightarrow Q) \rightarrow (Q \rightarrow R) \rightarrow P \rightarrow R$ trans- $\rightarrow f \ g \ p = g \ (f \ p)$ non-cont :  $(P \times (P \rightarrow \mathbb{O})) \rightarrow \mathbb{O}$ non-cont = elim- $\Sigma$  ( $\lambda \ p \ np \rightarrow np \ p$ )  $\mathsf{DeMorgan}_1: (P \to \mathbb{O}) + (Q \to \mathbb{O}) \to P \times Q \to \mathbb{O}$  $\mathsf{DeMorgan}_1 = \mathsf{elim}_+ (\lambda \ np \to \mathsf{elim}_{\Sigma} (\lambda \ p \ \_ \to np \ p))$  $(\lambda \ nq \rightarrow \text{elim-}\Sigma \ (\lambda \ q \ \rightarrow nq \ q))$  $\mathsf{DeMorgan}_2: (P \to \mathbf{0}) \times (Q \to \mathbf{0}) \to P + Q \to \mathbf{0}$  $DeMorgan_2 = elim-\Sigma (\lambda np nq \rightarrow elim-+ np nq)$  $\mathsf{DeMorgan}_3: (P + Q \to \mathbb{O}) \to (P \to \mathbb{O}) \times (Q \to \mathbb{O})$ **DeMorgan**<sub>3</sub>  $h = (\lambda \ p \rightarrow h \ (inl \ p))$ ,  $(\lambda \ q \rightarrow h \ (\text{inr} \ q))$ -- the final DeMorgan law cannot be proven constructively module Classical (excluded-middle :  $(S : Set) \rightarrow S + (S \rightarrow 0)$ ) where  $\mathsf{DeMorgan}_4: (P \times Q \to \mathbb{O}) \to (P \to \mathbb{O}) + (Q \to \mathbb{O})$  $\mathsf{DeMorgan}_4 h = \mathsf{elim}_+ (\lambda \ p \ \rightarrow \mathsf{elim}_+ \ (\lambda \ q \ \rightarrow \mathsf{elim}_- \mathbb{O} \ (h \ (p \ , q)))$  $(\lambda \ nq \rightarrow \text{inr} \ nq)$ (excluded-middle Q)) $(\lambda \ np \rightarrow \text{inl} \ np)$ (excluded-middle P)

```
{- PARAMETRICALLY-POLYMORPHIC FUNCTIONS -}
\mathsf{id}: \{A: \mathsf{Set}\; \alpha\} \to A \to A
id a = a
\_ : {A : Set \alpha} {B : A \to Set \beta} \rightarrow ((a : A) \rightarrow B a) \rightarrow (a : A) \rightarrow B a
_$_ = id
const : {A : Set \alpha} {B : Set \beta} \rightarrow A \rightarrow B \rightarrow A
const a \_ = a
if_then_else_: \{A : Set \ \alpha\} \rightarrow 2 \rightarrow A \rightarrow A \rightarrow A
if b then t else f = \text{elim} + (\text{const } t) (\text{const } f) b
flip : {A : Set \alpha} {B : Set \beta} {C : A \to B \to Set \gamma} \rightarrow
        ((a:A) \rightarrow (b:B) \rightarrow C \ a \ b) \rightarrow (b:B) \rightarrow (a:A) \rightarrow C \ a \ b)
flip f b a = f a b
\_\circ\_: \{A : \mathsf{Set} \ \alpha\} \ \{B : A \to \mathsf{Set} \ \beta\} \ \{C : \{a : A\} \to B \ a \to \mathsf{Set} \ \gamma\} \to
          (\{a:A\} \rightarrow (b:B a) \rightarrow C b) \rightarrow
          (g:(a:A) \rightarrow B a) \rightarrow
          (a:A) \rightarrow C (g a)
f \circ g = \lambda \ a \to f \ (g \ a)
\pi_1: {A : Set \alpha} {B : A \to Set \beta} \to \Sigma A B \to A
\pi_1 = \mathsf{elim} \cdot \Sigma \; (\lambda \; a \_ \rightarrow a)
\pi_2: {A : Set \alpha} {B : A \to Set \beta} \to (p : \Sigma A B) \to B (\pi_1 p)
\pi_2 = \operatorname{elim} \Sigma (\lambda \_ b \rightarrow b)
{- PROPOSITIONAL EQUALITY -}
data \_\equiv\_ \{A : \mathsf{Set} \ \alpha\} : A \to A \to \mathsf{Set} \ \alpha \text{ where}
   \mathsf{refl} = : (a : A) \to a \equiv a
```

```
\mathsf{elim} = \{A : \mathsf{Set} \ \alpha\} \ \{T : (a \ b : A) \to a \equiv b \to \mathsf{Set} \ \beta\} \ \{a \ b : A\} \to A
                ((c: A) \rightarrow T \ c \ c \ (\mathsf{refl} = c)) \rightarrow (p: a \equiv b) \rightarrow T \ a \ b \ p
\mathsf{elim} = f (\mathsf{refl} = a) = f a
sym-\equiv: {A : Set \alpha} {a \ b : A} \rightarrow a \equiv b \rightarrow b \equiv a
sym-\equiv = elim-\equiv (\lambda \ c \rightarrow refl-\equiv c)
trans-\equiv: {A : Set \alpha} {a \ b \ c : A} \rightarrow a \equiv b \rightarrow b \equiv c \rightarrow a \equiv c
trans-\equiv \{c = c\} = \text{elim} = (\lambda \ d \to \text{id} \ \{A = d \equiv c\})
\mathsf{cong}: \{A : \mathsf{Set} \ \alpha\} \ (P : A \to \mathsf{Set} \ \beta) \ \{a \ b : A\} \to P \ a \to a \equiv b \to P \ b
cong P = flip (elim= (\lambda \ c \rightarrow id \ \{A = P \ c\}))
\mathsf{resp} : \{A : \mathsf{Set} \ \alpha\} \ \{B : \mathsf{Set} \ \beta\} \ \{a \ b : A\} \rightarrow
           (f: A \to B) \to a \equiv b \to f \ a \equiv f \ b
resp f = \text{elim} = (\lambda \ c \rightarrow \text{refl} = (f \ c))
uniqueness : (u : 1) \rightarrow u \equiv \star
uniqueness = elim-1 (refl-\equiv \star)
Boolean : ((b: 2) \rightarrow (b \equiv \mathsf{true}) + (b \equiv \mathsf{false})) \times (\mathsf{true} \equiv \mathsf{false} \rightarrow \mathbb{0})
Boolean = true-or-false, absurd
    where
       true-or-false : (b: 2) \rightarrow (b \equiv \text{true}) + (b \equiv \text{false})
       true-or-false = elim-+ (elim-1 (inl (refl-\equiv true)))
                                                (elim-1 (inr (refl-\equiv false)))
       \mathsf{absurd} : \mathsf{true} \equiv \mathsf{false} \to \mathbb{0}
        absurd = cong (if_then 1 else 0) \star
{- PROPOSITIONAL EQUIVALENCE -}
\_\leftrightarrow\_: Set \alpha \to \mathsf{Set} \ \beta \to \mathsf{Set} \ (\alpha \sqcup \beta)
A \leftrightarrow B = (A \rightarrow B) \times (B \rightarrow A)
\mathsf{refl}{} \leftrightarrow : (A : \mathsf{Set} \ \alpha) \to A \leftrightarrow A
\operatorname{refl} \leftrightarrow A = \operatorname{id} , \operatorname{id}
```

$$\begin{split} & \text{sym-} \leftrightarrow : \{A : \text{Set } \alpha\} \{B : \text{Set } \beta\} \to A \leftrightarrow B \to B \leftrightarrow A \\ & \text{sym-} \leftrightarrow = \text{elim-} \Sigma (\lambda f g \to g, f) \\ & \text{trans-} \leftrightarrow : \{A : \text{Set } \alpha\} \{B : \text{Set } \beta\} \{C : \text{Set } \gamma\} \to \\ & A \leftrightarrow B \to B \leftrightarrow C \to A \leftrightarrow C \\ & \text{trans-} \leftrightarrow = \text{elim-} \Sigma (\lambda f g \to \text{elim-} \Sigma (\lambda f' g' \to f' \circ f, g \circ g')) \\ & \{- \text{ TYPE EQUIVALENCE } -\} \\ \underline{-} \cong : \text{Set } \alpha \to \text{Set } \beta \to \text{Set } (\alpha \sqcup \beta) \\ & A \simeq B = \Sigma (A \to B) (\lambda f \to \\ \Sigma (B \to A) (\lambda g \to ((a : A) \to a \equiv (g \circ f) a) \\ & \times ((b : B) \to b \equiv (f \circ g) b))) \\ & \text{equiv-3-3} = f, g, p, q \\ & \text{where} \\ & f : 1 + (1 + 1) \to (1 + 1) + 1 \\ & f = \text{elim-+} (\text{const (int (int *))}) \\ & (\text{const (int (int *)))} \\ & (\text{elim-+} (\text{elim-1 (refl-= (int (int *))))} \\ & (\text{elim-1 (refl-= (int (int *)))) \\ & (\text{elim-1 (refl-= (int *)))} \\ \end{aligned}$$

 $\begin{array}{l} \operatorname{refl-}\simeq:\,(A:\operatorname{Set}\,\alpha)\to A\simeq A\\ \operatorname{refl-}\simeq A=\operatorname{id}\,\text{, id , refl-}\equiv\,\text{, refl-}\equiv\end{array}$ 

 $\begin{array}{l} \mathsf{sym-} \simeq : \ \{A : \mathsf{Set} \ \alpha\} \ \{B : \mathsf{Set} \ \beta\} \to A \simeq B \to B \simeq A \\ \mathsf{sym-} \simeq = \mathsf{elim-} \mathbf{\Sigma} \ (\lambda \ f \to \mathsf{elim-} \mathbf{\Sigma} \ (\lambda \ g \to \mathsf{elim-} \mathbf{\Sigma} \ (\lambda \ p \ q \to g \ , f \ , q \ , p))) \end{array}$ 

$$\begin{array}{l} \operatorname{trans-} \simeq : \ \{A : \operatorname{Set} \alpha\} \ \{B : \operatorname{Set} \beta\} \ \{C : \operatorname{Set} \gamma\} \rightarrow \\ A \simeq B \rightarrow B \simeq C \rightarrow A \simeq C \\ \operatorname{trans-} \simeq (f \ , \ g \ , \ p \ , \ q) \ (f' \ , \ g' \ , \ p' \ , \ q') \\ = f' \circ f \ , \ g \circ g' \ , \ (\lambda \ a \rightarrow \operatorname{cong} \ (\lambda \ b \rightarrow a \equiv g \ b) \ (p \ a) \ (p' \ (f \ a))) \\ , \ (\lambda \ c \rightarrow \operatorname{cong} \ (\lambda \ b \rightarrow c \equiv f' \ b) \ (q' \ c) \ (q \ (g' \ c))) \end{array}$$

{- FUNCTION EQUIVALENCE -}

$$\_\sim\_: \{A : \mathsf{Set} \ \alpha\} \ \{B : A \to \mathsf{Set} \ \beta\} \to \\ (f \ g : (a : A) \to B \ a) \to \\ \mathsf{Set} \ (\alpha \sqcup \beta) \\ f \sim g = (a : \_) \to f \ a \equiv g \ a \end{cases}$$

$$\begin{array}{l} \mathsf{refl-}\sim:\;\{A:\mathsf{Set}\;\alpha\}\;\{B:A\to\mathsf{Set}\;\beta\}\to\\ (f:\,(a:\,A)\to B\;a)\to\\ f\sim f\\ \mathsf{refl-}\sim f\;a=\mathsf{refl-}{\equiv}\;(f\;a) \end{array}$$

$$\begin{array}{l} \mathsf{sym-}\sim:\ \{A:\mathsf{Set}\ \alpha\}\ \{B:A\to\mathsf{Set}\ \beta\}\to\\ &\ \{f\ g:\ (a:\ A)\to B\ a\}\to\\ &\ f\sim g\to g\sim f\\ \mathsf{sym-}\sim\ e\ a=\ \mathsf{sym-}\equiv\ (e\ a) \end{array}$$

 $\begin{array}{l} \operatorname{trans-} \sim : \ \{A : \operatorname{\mathsf{Set}} \alpha\} \ \{B : A \to \operatorname{\mathsf{Set}} \beta\} \to \\ \{f \ g \ h : \ (a : \ A) \to B \ a\} \to \\ f \sim g \to g \sim h \to f \sim h \\ \operatorname{trans-} \sim d \ e \ a = \operatorname{trans-} \equiv (d \ a) \ (e \ a) \end{array}$ 

 $\begin{array}{l} \mathsf{FunExt}: \ \{A:\mathsf{Set}\ \alpha\}\ \{B:A\to\mathsf{Set}\ \beta\}\to\\ (f\ g:\ (a:\ A)\to B\ a)\to\\ \mathsf{Set}\ (\alpha\sqcup\beta)\\ \end{array}$   $\begin{array}{l} \mathsf{FunExt}\ f\ g=f\sim g\to f\equiv g \end{array}$ 

{- INDUCTIVE TYPES -} data  $\mathbb{N}$  : Set where zero : ℕ  $\mathsf{succ}: \mathbb{N} \to \mathbb{N}$ elim- $\mathbb{N}$  : {  $T : \mathbb{N} \to \mathsf{Set} \ \alpha$  }  $\to T \ 0 \to$  $((n:\mathbb{N}) \rightarrow T \ n \rightarrow T \ (\mathsf{succ} \ n)) \rightarrow$  $(n:\mathbb{N}) \to T n$ elim- $\mathbb{N} z = 0$ = zelim- $\mathbb{N} z s$  (succ n) = s n (elim- $\mathbb{N} z s n$ ) inj-succ :  $\{m \ n : \mathbb{N}\} \rightarrow \text{succ} \ m \equiv \text{succ} \ n \rightarrow m \equiv n$ inj-succ  $\{m\} = \mathsf{cong} (\mathsf{elim} \cdot \mathbb{N} \ \mathbb{O} (\lambda \ o \_ \rightarrow m \equiv o)) (\mathsf{refl} \cdot \equiv m)$ pred :  $\mathbb{N} \to \mathbb{N}$ pred = elim- $\mathbb{N}$  zero ( $\lambda \ n \_ \rightarrow n$ )  $\mathsf{add}\,:\,\mathbb{N}\to\mathbb{N}\to\mathbb{N}$ add  $n = \operatorname{elim} \mathbb{N} n \ (\lambda \_ \to \operatorname{succ})$ add-id-r : flip add 0  $\sim$  id  $add-id-r = refl-\equiv$ add-id-I : add 0  $\sim$  id add-id-I = elim- $\mathbb{N}$  (refl- $\equiv$  0) ( $\lambda \_ \rightarrow$  resp succ) data List (A : Set) : Set where [] : List *A*  $\_::\_: A \rightarrow \mathsf{List} A \rightarrow \mathsf{List} A$ data Fin :  $\mathbb{N} \rightarrow \mathsf{Set}$  where Fin-zero :  $\{n : \mathbb{N}\} \rightarrow \mathsf{Fin} (\mathsf{succ} n)$ Fin-succ :  $\{n : \mathbb{N}\} \rightarrow \text{Fin } n \rightarrow \text{Fin} (\text{succ } n)$ data  $\_\leq\_: \mathbb{N} \to \mathbb{N} \to \mathsf{Set}$  where  $\operatorname{zero-} < : \{n : \mathbb{N}\} \to \operatorname{zero} < n$  $\operatorname{succ} \leq \{m \ n : \mathbb{N}\} \to m \leq n \to \operatorname{succ} m \leq \operatorname{succ} n$ 

```
\mathsf{half}:\,\mathbb{N}\to\mathbb{N}
half zero
                = zero
half (succ zero) = zero
half (succ (succ n)) = succ (half n)
\mathsf{refl} \le : (n : \mathbb{N}) \to n \le n
refl \le zero = zero \le
refl-\leq (succ n) = succ-\leq (refl-\leq n)
\mathsf{trans-} \leq : \ \{m \ n \ o : \ \mathbb{N}\} \rightarrow m \leq n \rightarrow n \leq o \rightarrow m \leq o
trans-\leq zero-\leq q = zero-\leq
trans-\leq (succ-\leq p) (succ-\leq q) = succ-\leq (trans-\leq p q)
{- COINDUCTIVE TYPES -}
record Stream (A : Set) : Set where
  coinductive
  field
     head : A
     tail : Stream A
from : \mathbb{N} \to \mathsf{Stream} \ \mathbb{N}
head (from n) = n
tail (from n) = from (succ n)
nats : Stream ℕ
nats = from 0
-- cannot be defined for streams
-- due to the unguarded recursive call
filter : \{A : \mathsf{Set}\} \to (A \to 2) \to \mathsf{List} \ A \to \mathsf{List} \ A
filter f[] = []
filter f(h :: t) = if f h then h :: filter f t else filter f t
\mathsf{map}: \{A \ B : \mathsf{Set}\} \to (A \to B) \to \mathsf{Stream} \ A \to \mathsf{Stream} \ B
head (map f s) = f (head s)
tail (map f s) = map f (tail s)
```

```
zip: \{A \ B : Set\} \rightarrow Stream \ A \rightarrow Stream \ B \rightarrow Stream \ (A \times B)
head (zip \ a \ b) = head \ a, head b
tail (zip \ a \ b) = zip (tail \ a) (tail \ b)
record Always \{A : Set\} (P : A \rightarrow Set) (s : Stream A) : Set where
   coinductive
   field
     always-head : P (head s)
     always-tail : Always P (tail s)
record Bisim \{A : Set\} (a \ b : Stream \ A) : Set where
   coinductive
   field
     bisim-head : head a \equiv head b
     bisim-tail : Bisim (tail a) (tail b)
map-succ-from : (n : \mathbb{N}) \to \text{Bisim} (\text{map succ} (\text{from } n)) (tail (from n))
bisim-head (map-succ-from n) = refl-\equiv
                                                            (succ n)
bisim-tail (map-succ-from n) = map-succ-from (succ n)
refl-Bisim : \{A : Set\} (a : Stream A) \rightarrow Bisim a a
bisim-head (refl-Bisim a) = refl-\equiv
                                              (head a)
bisim-tail (refl-Bisim a) = refl-Bisim (tail a)
sym-Bisim : \{A : \mathsf{Set}\} \{a \ b : \mathsf{Stream} \ A\} \to \mathsf{Bisim} \ a \ b \to \mathsf{Bisim} \ b \ a
bisim-head (sym-Bisim p) = sym-\equiv
                                                 (bisim-head p)
bisim-tail (sym-Bisim p) = sym-Bisim (bisim-tail p)
trans-Bisim : \{A : \mathsf{Set}\} \{a \ b \ c : \mathsf{Stream} \ A\} \rightarrow
                 \mathsf{Bisim} \ a \ b \to \mathsf{Bisim} \ b \ c \to \mathsf{Bisim} \ a \ c
bisim-head (trans-Bisim p q) = trans-\equiv
                                                      (bisim-head p)
                                                      (bisim-head q)
```

bisim-tail (trans-Bisim p q) = trans-Bisim (bisim-tail p)

```
(bisim-tail q)
```

## A.2 Epistemic Modal Logic

```
{-
   CHAPTER III: EPISTEMIC MODAL LOGIC
   - checked with Agda version 2.7.0.1
-}
```

open import Agda.Primitive open import TypeTheory

module EpistemicLogic ( $\sigma$  : Level) (*State* : Set  $\sigma$ ) where

{- EVENTS AND EPISTEMIC LOGIC CONNECTIVES -}

Event : Set  $(\sigma \sqcup \mathsf{Isuc} \alpha)$ Event  $\{\alpha\} = State \to \mathsf{Set} \alpha$ 

 $\perp$  : Event  $\perp$  = const 0

 $\top$ : Event  $\top$  = const 1

 $\neg\_: \mathsf{Event} \{\alpha\} \to \mathsf{Event}$  $\neg \phi = \lambda \ w \to \phi \ w \to \mathbf{0}$ 

 $\_\lor\_: \mathsf{Event} \ \{\alpha\} \to \mathsf{Event} \ \{\beta\} \to \mathsf{Event} \ \\ \phi \lor \psi = \lambda \ w \to \phi \ w + \psi \ w$ 

 $\_\wedge\_: \mathsf{Event} \{\alpha\} \to \mathsf{Event} \{\beta\} \to \mathsf{Event}$  $\phi \land \psi = \lambda \ w \to \phi \ w \times \psi \ w$ 

\_⇒\_: Event {α} → Event {β} → Event  $\phi \Rightarrow \psi = \lambda \ w \rightarrow \phi \ w \rightarrow \psi \ w$ 

\_⇔\_ : Event {α} → Event {β} → Event  $\phi \Leftrightarrow \psi = \lambda \ w \to \phi \ w \leftrightarrow \psi \ w$   $\begin{array}{l} \forall \forall \_ : \; \mathsf{Event} \; \{\alpha\} \to \mathsf{Set} \; (\sigma \sqcup \alpha) \\ \forall \forall \_ \phi = \{w : \; State\} \to \phi \; w \\ \\ \_ \subseteq \_ : \; \mathsf{Event} \; \{\alpha\} \to \mathsf{Event} \; \{\beta\} \to \mathsf{Set} \; (\sigma \sqcup \alpha \sqcup \beta) \\ \phi \subseteq \psi = \forall \forall \; (\phi \Rightarrow \psi) \\ \\ \_ \approx \_ : \; \mathsf{Event} \; \{\alpha\} \to \mathsf{Event} \; \{\beta\} \to \mathsf{Set} \; (\sigma \sqcup \alpha \sqcup \beta) \\ \phi \approx \psi = \forall \forall \; (\phi \Leftrightarrow \psi) \end{array}$ 

{- KNOWLEDGE OPERATOR SEMANTICS -}

 $\_\models\_: \{A : \mathsf{Set} \ \alpha\} \to (A \to \mathsf{Event} \ \{\beta\}) \to \mathsf{Event} \ \{\gamma\} \to \mathsf{Set} \ (\sigma \sqcup \alpha \sqcup \beta \sqcup \gamma) \\ \varPhi \models \psi = \{w : \ State\} \to ((a : \_) \to \varPhi \ a \ w) \to \psi \ w$ 

```
record KOp : Set (\sigma \sqcup lsuc (\alpha \sqcup \beta)) where

field

K : Event \{\alpha\} \rightarrow Event \{\alpha\}

preserves-\models : \{A : Set \ \beta\} \ (\Phi : A \rightarrow Event) \rightarrow

\{\psi : Event\} \rightarrow \Phi \models \psi \rightarrow \mathsf{K} \circ \Phi \models \mathsf{K} \psi

axiom-T : \{\phi : Event\} \rightarrow \ \mathsf{K} \phi \subseteq \phi

axiom-4 : \{\phi : Event\} \rightarrow \ \mathsf{K} \phi \subseteq \mathsf{K} \ (\mathsf{K} \phi)

axiom-5 : \{\phi : Event\} \rightarrow \neg \ \mathsf{K} \phi \subseteq \mathsf{K} \ (\neg \ \mathsf{K} \phi)
```

 $\begin{array}{l} \mathsf{K}\text{-gen} : \{\{\_: \mathsf{K}\mathsf{Op}\}\} \to \{\phi : \mathsf{Event} \ \{\alpha\}\} \to \forall \forall \ \phi \to \forall \forall \ \mathsf{K} \ \phi \\ \\ \mathsf{K}\text{-gen} \ h = \mathsf{preserves}\text{-}\vDash \mathsf{elim}\text{-}\mathbb{0} \ (\mathsf{const} \ h) \ \mathsf{elim}\text{-}\mathbb{0} \end{array}$ 

axiom-K : {{\_: KOp}}  $\rightarrow \{\phi \ \psi : \text{Event } \{\alpha\}\} \rightarrow \mathsf{K} \ (\phi \Rightarrow \psi) \subseteq (\mathsf{K} \ \phi \Rightarrow \mathsf{K} \ \psi)$ axiom-K {\_} { $\phi$ } { $\psi$ }  $p \ q = \mathsf{preserves} \models (\mathsf{if\_then} \ \phi \Rightarrow \psi \ \mathsf{else} \ \phi)$  $(\lambda \ f \rightarrow f \ \mathsf{true} \ (f \ \mathsf{false}))$  $(\mathsf{elim} + (\mathsf{const} \ p) \ (\mathsf{const} \ q))$ 

{- CORRESPONDENCE WITH RELATIONAL SEMANTICS -}

Kripke :  $\{\alpha : \text{Level}\} \rightarrow \text{Set} (\sigma \sqcup \text{Isuc } \alpha)$ Kripke  $\{\alpha\} = \text{EquivRel} \{\sigma\} \{\alpha\} State$ 

 $\mathsf{K}[\_]: (State \to State \to \mathsf{Set} \ \alpha) \to \mathsf{Event} \ \{\beta\} \to \mathsf{Event} \ \{\sigma \sqcup \alpha \sqcup \beta\}$  $\mathsf{K}[\_R\_] \ \phi \ w = \{v: \ State\} \to w \ R \ v \to \phi \ v$ 

 $R[\_]_: State \to (\mathsf{Event} \{\alpha\} \to \mathsf{Event} \{\beta\}) \to State \to \mathsf{Set} (\sigma \sqcup \mathsf{Isuc} \ \alpha \sqcup \beta)$  $w \ \mathsf{R}[K] \ v = \{\phi : \mathsf{Event}\} \to K \ \phi \ w \to K \ \phi \ v$ 

```
\mathsf{Kripke}\mathsf{-}\mathsf{KOp}: \{\{\_: \mathsf{Kripke} \ \{\alpha\}\}\} \to \mathsf{KOp} \ \{\sigma \sqcup \alpha\} \ \{\beta\}
Kripke-KOp = record \{ K \}
                                                           = K[\_R_]
                                        preserves-\models = \lambda \Phi h p q \rightarrow h (\lambda a \rightarrow p a q)
                                       axiom-T = \lambda k \rightarrow k refl
                                        axiom-4 = \lambda k p \rightarrow k \circ \text{trans } p
                                        axiom-5 = \lambda h p k \rightarrow h (k \circ \text{trans (sym } p))
module Classical (elim-\neg \neg: {A : Set \alpha} \rightarrow ((A \rightarrow 0) \rightarrow 0) \rightarrow A) where
   KOp-Kripke : {{ : KOp {\alpha} {\beta}}} \rightarrow Kripke
   KOp-Kripke = record { R_ = R[K]
                                                                                                                  ;
                                           refl = id
                                                                                                                  ÷
                                           sym = \lambda p k \rightarrow elim \neg \neg (flip axiom-T k
                                                                                        \circ p \circ axiom-5);
                                           trans = \lambda p q \rightarrow q \circ p
                                                                                                                  }
iso-R : {\beta : Level} {{\_ : Kripke {\sigma \sqcup \alpha}} (w v : State) \rightarrow
            w \mathrel{\mathsf{R}} v \leftrightarrow w \mathrel{\mathsf{R}} [ \mathrel{\mathsf{K}} [ \_ \mathrel{\mathsf{R}} \_ ] ] v
iso-R {\beta} w v = ltr , rtl
   where
      \mathsf{Itr}: w \mathsf{R} v \to w \mathsf{R}[\mathsf{K}[\mathsf{R}]] v
      Itr p \ k = axiom-4 {{Kripke-KOp {\beta = \beta}} k \ p
       \mathsf{rtl}: w \mathsf{R}[\mathsf{K}[\_\mathsf{R}\_]] v \to w \mathsf{R} v
       rtl p = p id refl
iso-K : {{_ : KOp}} (\phi : Event {\alpha}) \rightarrow K \phi \approx K[ _R[ K ]_ ] \phi
iso-K \phi = \mathsf{Itr} , rtl
   where
      \mathsf{ltr}: \mathsf{K} \phi \subseteq \mathsf{K}[\_\mathsf{R}[\mathsf{K}]\_] \phi
      Itr k p = axiom-T (p k)
      \mathsf{rtl}:\mathsf{K}[\_\mathsf{R}[\mathsf{K}]\_]\phi\subseteq\mathsf{K}\phi
       rtl k = preserves-\models (K \circ \pi_1)
                                       (\lambda \ h \rightarrow k \ (\lambda \ \{\psi\} \ j \rightarrow h \ (\psi \ , \ j)))
                                       (axiom-4 \circ \pi_2)
```

# A.3 Coinductive Common Knowledge

```
{-
   CHAPTER IV: COINDUCTIVE COMMON KNOWLEDGE
   - checked with Agda version 2.7.0.1
-}
```

open import Agda.Primitive open import TypeTheory

module CommonKnowledge (State : Set) (Agent : Set) (agent : Agent) (K-rel : Agent  $\rightarrow$  EquivRel { $\beta =$  Izero} State) where

open import EpistemicLogic Izero State

 $\_\propto[\_]\_: State \to Agent \to State \to Set$  $w \propto [a] v = (K-rel a).\_R\_w v$ 

 $\mathsf{K} : Agent \to \mathsf{Event} \{\alpha\} \to \mathsf{Event}$  $\mathsf{K} \ a = \mathsf{K}[\_\infty[ \ a \ ]\_ ]$ 

 $\begin{array}{l} \mathsf{K}\text{-op} : Agent \to \mathsf{KOp} \\ \mathsf{K}\text{-op} \ a = \mathsf{Kripke}\text{-}\mathsf{KOp} \ \{\mathsf{Izero}\} \ \{\mathsf{Izero}\} \ \{\mathsf{K}\text{-}\mathit{rel} \ a\}\} \end{array}$ 

 $\begin{array}{l} \mathsf{EK} : \mathsf{Event} \ \{\alpha\} \to \mathsf{Event} \\ \mathsf{EK} \ \phi \ w = \{a: \mathit{Agent}\} \to \mathsf{K} \ a \ \phi \ w \end{array}$ 

axiom-T-EK :  $\{\phi : \text{Event}\} \rightarrow \text{EK } \phi \subseteq \phi$ axiom-T-EK  $k = \text{axiom-T} \{\{\text{K-op } agent\}\} k$ 

```
{- COINDUCTIVE COMMON KNOWLEDGE DEFINITION -}
record cCK (\phi : Event) (w : State) : Set where
   coinductive
   field
     \mathsf{cCK}-\mathsf{EK} : \mathsf{EK} \phi w
     cCK-cCKEK : cCK (EK \phi) w
cCK-EKcCK : {\phi : Event} \rightarrow cCK \phi \subseteq EK (cCK \phi)
                  (cCK-EKcCK \ k \ p) = cCK-EK (cCK-cCKEK \ k) \ p
cCK-EK
cCK-cCKEK (cCK-EKcCK k p) = cCK-EKcCK (cCK-cCKEK k) p
\operatorname{recCK} : Event \{\alpha\} \to \mathbb{N} \to \operatorname{Event}
\operatorname{rec}\mathsf{CK}\phi zero =\phi
\operatorname{recCK} \phi (\operatorname{succ} n) = \operatorname{recCK} (\operatorname{EK} \phi) n
\operatorname{recCK-cCK} : (\phi : Event) \rightarrow \operatorname{recCK} \phi \models \operatorname{cCK} \phi
                  (\operatorname{rec}\mathsf{CK}\operatorname{-c}\mathsf{CK}\phi h) = h\ 1
cCK-EK
cCK-cCKEK (recCK-cCK \phi h) = recCK-cCK (EK \phi) (h \circ succ)
\mathsf{cCK}\operatorname{-recCK}: (\phi : \mathsf{Event}) \to (n : \mathbb{N}) \to \mathsf{cCK} \ \phi \subseteq \mathsf{recCK} \ \phi \ n
cCK-recCK \phi zero k = axiom-T-EK (cCK-EK k)
cCK-recCK \phi (succ n) k = cCK-recCK (EK \phi) n (cCK-cCKEK k)
{- RELATIONAL COMMON KNOWLEDGE DEFINITION -}
data \_\infty\_: State \rightarrow State \rightarrow Set where
   union-\infty: {a : Agent} {w v : State} \rightarrow w \propto [a] v \rightarrow w \propto v
  \mathsf{trans-} \infty : \{ w \ v \ u : \ State \} \to w \propto v \to v \propto u \to w \propto u
\mathsf{rCK} : \mathsf{Event} \{\alpha\} \to \mathsf{Event}
rCK = K[\_\infty\_]
\mathsf{refl}-\infty : {w : State} \rightarrow w \propto w
refl-\infty = union-\infty ((K-rel agent), refl)
\mathsf{sym}\text{-}\infty: \{w \ v: \ State\} \to w \propto v \to v \propto w
sym-\propto (union-\propto {a} p) = union-\propto ((K-rel a).sym p)
sym-\infty (trans-\infty p q) = trans-\infty (sym-\infty q) (sym-\propto p)
```

```
Kripke-\infty : Kripke
Kripke-\propto = record { _R_ = _\propto_ ;
                                refl = refl-\propto ;
                                sym = sym-\propto ;
                                trans = trans-\infty }
rCK-op : KOp
rCK-op = Kripke-KOp {\beta = lzero} {{Kripke-\infty}}
\mathsf{rCK}-EK : {\phi : Event {\alpha}} \rightarrow \mathsf{rCK} \ \phi \subseteq \mathsf{EK} \ \phi
rCK-EK r p = r (union-\propto p)
\mathsf{rCK}-\mathsf{rCKEK} : \{\phi : \mathsf{Event} \{\alpha\}\} \to \mathsf{rCK} \phi \subseteq \mathsf{rCK} (\mathsf{EK} \phi)
rCK-rCKEK r \ p \ q = r \ (\text{trans-} \propto p \ (\text{union-} \propto q))
\mathsf{rCK}-\mathsf{EKrCK} : {\phi : Event {\alpha}} \rightarrow \mathsf{rCK} \phi \subseteq \mathsf{EK} (\mathsf{rCK} \phi)
rCK-EKrCK r p q = r (trans-\propto (union-\propto p) q)
{- EQUIVALENCE WITH RELATIONAL DEFINITION -}
\mathsf{cCK}\text{-transport}: \{\phi : \mathsf{Event}\} \{w \ v : \ State\} \to \mathsf{cCK} \ \phi \ w \to w \propto v \to \mathsf{cCK} \ \phi \ v
cCK-transport k (union-\propto \{a\} p) = axiom-T {{K-op a}}
                                                  axiom-4 \{\{K-op a\}\} (cCK-EKcCK k) p
cCK-transport k (trans-\propto p q) = cCK-transport (cCK-transport k p) q
\mathsf{cCK}-\mathsf{rCK} : {\phi : Event} \rightarrow \mathsf{cCK} \ \phi \approx \mathsf{rCK} \ \phi
cCK-rCK = ltr , rtl
   where
      \mathsf{ltr}: \{\phi: \mathsf{Event}\} \to \mathsf{cCK} \ \phi \subseteq \mathsf{rCK} \ \phi
      Itr k p = axiom-T-EK (cCK-EK (cCK-transport k p))
      \mathsf{rtl}: \{\phi: \mathsf{Event}\} \to \mathsf{rCK} \ \phi \subset \mathsf{cCK} \ \phi
      cCK-EK
                         (\mathsf{rtl} \ k) = \mathsf{rCK}\mathsf{-}\mathsf{EK} \ k
      cCK-cCKEK (rtl k) = rtl (rCK-rCKEK k)
```
{- COINDUCTIVE COMMON KNOWLEDGE IS A KNOWLEDGE OPERATOR -}

preserves-
$$\models$$
-cCK : { $A$  : Set  $\alpha$ } ( $\Phi$  :  $A \rightarrow \text{Event}$ ) { $\psi$  : Event}  $\rightarrow \Phi \models \psi \rightarrow \text{cCK} \circ \Phi \models \text{cCK} \psi$   
preserves- $\models$ -cCK  $\Phi$   $h$   $i = \pi_2$  cCK-rCK ( $\lambda \ p \rightarrow h$  ( $\lambda \ t \rightarrow \text{axiom-T}$  {{rCK-op}}  
 $\circ \pi_1$  cCK-rCK  
\$ cCK-transport ( $i \ t$ )  $p$ ))

```
axiom-T-cCK : {\phi : Event} \rightarrow cCK \phi \subseteq \phi
axiom-T-cCK k = axiom-T {{rCK-op}} (\pi_1 cCK-rCK k)
```

```
\begin{array}{l} \operatorname{axiom-4-cCK} : \{\phi : \operatorname{Event}\} \to \operatorname{cCK} \phi \subseteq \operatorname{cCK} (\operatorname{cCK} \phi) \\ \operatorname{axiom-4-cCK} k = \pi_2 \operatorname{cCK-rCK} (\lambda \ p \to \pi_2 \operatorname{cCK-rCK} \\ (\lambda \ q \to \operatorname{axiom-T-cCK} \\ \$ \quad \operatorname{cCK-transport} (\operatorname{cCK-transport} k \ p) \ q)) \end{array}
```

```
\begin{array}{l} \operatorname{axiom-5-cCK} : \{\phi : \operatorname{Event}\} \to \neg \operatorname{cCK} \phi \subseteq \operatorname{cCK} (\neg \operatorname{cCK} \phi) \\ \operatorname{axiom-5-cCK} nk = \pi_2 \operatorname{cCK-rCK} (\lambda \ p \ k \to nk \\ \circ \ \operatorname{cCK-transport} k \\ \$ \ \operatorname{sym-x} p) \end{array}
```

```
cCK-op : KOp \{Izero\} \{\beta\}
cCK-op = record \{ K = cCK ;
preserves-\models = preserves-\models-cCK ;
axiom-T = axiom-T-cCK ;
axiom-4 = axiom-4-cCK ;
axiom-5 = axiom-5-cCK \}
```

## Bibliography

- Thomas Ågotnes and Natasha Alechina. Coalition logic with individual, distributed and common knowledge. *Journal of Logic and Computation*, 29:1041–1069, 2018.
- [2] Robert J. Aumann. Agreeing to disagree. Annals of Statistics, 4:1236–1239, 1976.
- [3] Robert J. Aumann. Backward induction and common knowledge of rationality. Games and Economic Behavior, 8(1):6–19, 1995.
- [4] Jon Barwise. Three views of common knowledge. In Proceedings of the 2nd Conference on Theoretical Aspects of Reasoning about Knowledge, Pacific Grove, CA, March 1988, pages 365–379, 1988.
- [5] Colm Baston and Venanzio Capretta. The coinductive formulation of common knowledge. In International Conference on Interactive Theorem Proving, 2018.
- [6] Colm Baston and Venanzio Capretta. Game forms for coalition effectivity functions. In International Conference on Types for Proofs and Programs, 2019.
- [7] Yves Bertot and Pierre Castran. Interactive theorem proving and program development: Coq'art the calculus of inductive constructions. 2010.
- [8] Mark Bickford, Fedor A. Bogomolov, and Yuri Tschinkel. Vladimir Voevodsky – work and destiny. 2017.
- [9] Edwin C. Brady. Type-driven development with Idris. Manning, 2017.
- [10] John Charles Burkill and John Edensor Littlewood. A mathematician's miscellany. *The Mathematical Gazette*, 38:47, 1954.
- [11] Kevin Buzzard. What is the Xena project? https://xenaproject. wordpress.com/what-is-the-xena-project/.

- [12] Venanzio Capretta. Common knowledge as a coinductive modality. In Erik Barendsen, Herman Geuvers, Venanzio Capretta, and Milad Niqui, editors, *Reflections on Type Theory, Lambda Calculus, and the Mind*, pages 51–61. ICIS, Faculty of Science, Radbout University Nijmegen, 2007. Essays Dedicated to Henk Barendregt on the Occasion of his 60th Birthday.
- [13] Venanzio Capretta. Coalgebras in functional programming and type theory. *Theoretical Computer Science*, 412(38):5006–5024, 2011. CMCS Tenth Anniversary Meeting.
- [14] Alonzo Church. A set of postulates for the foundation of logic. Annals of Mathematics, 33:346, 1932.
- [15] Alonzo Church. An unsolvable problem of elementary number theory. American Journal of Mathematics, 58:345, 1936.
- [16] Alonzo Church. A formulation of the simple theory of types. Journal of Symbolic Logic, 5:56 – 68, 1940.
- [17] Cyril Cohen, Thierry Coquand, Simon Huber, and Anders Mörtberg. Cubical type theory: a constructive interpretation of the univalence axiom. 2016.
- [18] Thierry Coquand. Pattern matching with dependent types. 1992.
- [19] Thierry Coquand. Infinite objects in type theory. In Henk Barendregt and Tobias Nipkow, editors, *Types for Proofs and Programs. International Work*shop TYPES'93, volume 806 of Lecture Notes in Computer Science, pages 62–78. Springer-Verlag, 1993.
- [20] Paolo Crivelli, Timothy Luke Williamson, Gareth Hughes, and Max J. Cresswell. A new introduction to modal logic. 1998.
- [21] Haskell B. Curry. Some philosophical aspects of combinatory logic. Studies in logic and the foundations of mathematics, 101:85–101, 1980.
- [22] Leonardo Mendonça de Moura and Sebastian Ullrich. The Lean 4 theorem prover and programming language. In *CADE*, 2021.
- [23] Peter Dybjer. Inductive sets and families in Martin-Löf's type theory and their set-theoretic semantics. 1991.
- [24] Ronald Fagin, Joseph Y. Halpern, Moshe Y. Vardi, and Yoram Moses. *Reasoning About Knowledge*. MIT Press, Cambridge, MA, USA, 1995.

- [25] Daniel P. Friedman and David S. Wise. Cons should not evaluate its arguments. In International Colloquium on Automata, Languages and Programming, 1976.
- [26] George Gamow and Marvin Stern. Puzzle Math. Viking Press, New York, 1958.
- [27] James Garson. Modal Logic. In Edward N. Zalta and Uri Nodelman, editors, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Spring 2024 edition, 2024.
- [28] Eduardo Giménez. Codifying guarded definitions with recursive schemes. In Peter Dybjer, Bengt Nordström, and Jan Smith, editors, *Types for Proofs* and Programs. International Workshop TYPES '94, volume 996 of Lecture Notes in Computer Science, pages 39–59. Springer, 1994.
- [29] Jean-Yves Girard. Une extension de Einterpretation de Gödel a Eanalyse, et son application a Eelimination des coupures dans Eanalyse et la theorie des types. Studies in logic and the foundations of mathematics, 63:63–92, 1971.
- [30] Jean-Yves Girard. Interpretation fonctionelle et elimination des coupures dans l'aritmetique d'ordre superieur. 1972.
- [31] Healfdene Goguen, Conor McBride, and James McKinna. Eliminating dependent pattern matching. In *Essays Dedicated to Joseph A. Goguen*, 2006.
- [32] Peter B. Henderson and James H. Morris. A lazy evaluator. Proceedings of the 3rd ACM SIGACT-SIGPLAN symposium on Principles on programming languages, 1976.
- [33] Jaakko Hintikka. Knowledge and Belief. Ithaca: Cornell University Press, 1962.
- [34] Ralf Hinze. The Bird tree. Journal of Functional Programming, 19:491 508, 2009.
- [35] Simon Huber. Canonicity for cubical type theory. Journal of Automated Reasining, pages 173–210, 2019.
- [36] Antonius J. C. Hurkens. A simplification of Girard's paradox. In International Conference on Typed Lambda Calculus and Applications, 1995.
- [37] Saul A. Kripke. A completeness theorem in modal logic. Journal of Symbolic Logic, 24(1):1–14, 03 1959.

- [38] Leslie Lamport, Robert E. Shostak, and Marshall C. Pease. The byzantine generals problem. ACM Trans. Program. Lang. Syst., 4:382–401, 1982.
- [39] Pierre Lescanne. Common knowledge logic in a higher order proof assistant. In Andrei Voronkov and Christoph Weidenbach, editors, *Programming Logics* - *Essays in Memory of Harald Ganzinger*, volume 7797 of *Lecture Notes in Computer Science*, pages 271–284. Springer, 2013.
- [40] C.I. Lewis and C.H. Langford. Symbolic Logic. Century philosophy series. Century Company, 1932.
- [41] David Lewis. Convention: A philosophical study. 1969.
- [42] Simon Marlow. Haskell 2010 language report. 2010.
- [43] Per Martin-Löf. An intuitionistic theory of types. 1972.
- [44] Per Martin-Löf. An intuitionistic theory of types: Predicative part. Studies in logic and the foundations of mathematics, 80:73–118, 1975.
- [45] Per Martin-Löf. Intuitionistic type theory. In Studies in proof theory, 1984.
- [46] Randall Munroe. XKCD: Blue Eyes. https://xkcd.com/blue\_eyes.html.
- [47] Hiroshi Nakano. A modality for recursion. Proceedings Fifteenth Annual IEEE Symposium on Logic in Computer Science (Cat. No.99CB36332), pages 255–266, 2000.
- [48] Ulf Norell. Dependently typed programming in Agda. In Lecture Notes from the Summer School on Advanced Functional Programming, AFP 2008, pages 230–266. Springer Berlin Heidelberg, 2009.
- [49] Marc Pauly. A modal logic for coalitional power in games. Journal of Logic and Computation, 12, 02 2002.
- [50] John C. Reynolds. Towards a theory of type structure. In Symposium on Programming, 1974.
- [51] John C. Reynolds. User-Defined Types and Procedural Data Structures as Complementary Approaches to Data Abstraction, pages 309–317. Springer-Verlag, New York, NY, 1978.
- [52] John C. Reynolds. Types, abstraction and parametric polymorphism. In *IFIP Congress*, 1983.

- [53] Bertrand Russell and Alfred North Whitehead. Principia Mathematica, volume 1. 1910.
- [54] Thomas Streicher. Semantics of type theory correctness, completeness and independence results. In *Progress in theoretical computer science*, 1991.
- [55] The Univalent Foundations Program. Homotopy Type Theory: Univalent Foundations of Mathematics. https://homotopytypetheory.org/book, Institute for Advanced Study, 2013.
- [56] Wiebe van der Hoek and Michael Wooldridge. Cooperation, knowledge, and time: Alternating-time temporal epistemic logic and its applications. *Studia Logica*, 75:125–157, 2003.
- [57] Philip Wadler. Theorems for free! In Conference on Functional Programming Languages and Computer Architecture, 1989.