

# A General Framework for Relational Parametricity

Kristina Sojakova  
Appalachian State University  
Boone, NC, USA  
sojakovak@appstate.edu

Patricia Johann  
Appalachian State University  
Boone, NC, USA  
johannp@appstate.edu

## Abstract

Reynolds' original theory of *relational parametricity* was intended to capture the idea that polymorphically typed System F programs preserve all relations between inputs. But as Reynolds himself later showed, his theory can only be formalized in a meta-theory with an impredicative universe, such as the Calculus of Inductive Constructions. Abstracting from Reynolds' ideas, Dunphy and Reddy developed their well-known framework for parametricity that uses parametric limits in reflexive graph categories and aims to subsume a variety of parametric models. As we observe, however, their theory is not sufficiently general to subsume the very model that inspired parametricity, namely Reynolds' original model, expressed inside type theory.

To correct this, we develop an abstract framework for relational parametricity that generalizes the notion of a reflexive graph categories and delivers Reynolds' model as a direct instance in a natural way. This framework is uniform with respect to a choice of meta-theory, which allows us to obtain the well-known PER model of Longo and Moggi as a direct instance in a natural way as well. In addition, we offer two novel relationally parametric models of System F: *i) a categorical version of Reynolds' model*, where types are functorial on isomorphisms and all polymorphic functions respect the functorial action, and *ii) a proof-relevant categorical version of Reynolds' model* (after Orsanigo), where, additionally, witnesses of relatedness are themselves suitably related. We show that, unlike previously existing frameworks for parametricity, ours recognizes both of these new models in a natural way. Our framework is thus *descriptive*, in that it accounts for well-known models, as well as *prescriptive*, in that it identifies abstract properties that good models of relational parametricity should satisfy and suggests new constructions of such models.

**CCS Concepts** • Theory of computation → Type theory:

**Keywords** System F, categorical semantics, parametricity

## ACM Reference format:

Kristina Sojakova and Patricia Johann. 2018. A General Framework for Relational Parametricity. In *Proceedings of LICS '18: 33rd Annual ACM/IEEE Symposium on Logic in Computer Science, Oxford, United Kingdom, July 9–12, 2018 (LICS '18)*, 10 pages. <https://doi.org/10.1145/3209108.3209141>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

LICS '18, July 9–12, 2018, Oxford, United Kingdom

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-5583-4/18/07...\$15.00

<https://doi.org/10.1145/3209108.3209141>

## 1 Introduction

Reynolds [13] introduced the notion of *relational parametricity* to model the extensional behavior of programs in System F [6], the formal calculus at the core of all polymorphic functional languages. His goal was to give a type  $\alpha \vdash T(\alpha)$  an *object interpretation*  $T_0$  and a *relational interpretation*  $T_1$ , where  $T_0$  takes sets to sets and  $T_1$  takes relations  $R \subseteq A \times B$  to relations  $T_1(R) \subseteq T_0(A) \times T_0(B)$ . A term  $\alpha; x : S(\alpha) \vdash t(\alpha, x) : T(\alpha)$  was to be interpreted as a map  $t_0$  associating to each set  $A$  a function  $t_0(A) : S_0(A) \rightarrow T_0(A)$ . The interpretations were to be given inductively on the structure of  $T$  and  $t$  in such a way that they implied two key theorems: the *Identity Extension Lemma*, stating that if  $R$  is the equality relation on  $A$  then  $T_1(R)$  is the equality relation on  $T_0(A)$ ; and the *Abstraction Theorem*, stating that, for any relation  $R \subseteq A \times B$ ,  $t_0(A)$  and  $t_0(B)$  map arguments related by  $S_1(R)$  to results related by  $T_1(R)$ . A similar result holds for types and terms with any number of free variables.

In Reynolds' treatment of relational parametricity, if  $U(\alpha)$  is the type  $\alpha \vdash S(\alpha) \rightarrow T(\alpha)$ , for example, then  $U_0(A)$  is the set of functions  $f : S_0(A) \rightarrow T_0(A)$  and, for  $R \subseteq A \times B$ ,  $U_1(R)$  relates  $f : S_0(A) \rightarrow T_0(A)$  to  $g : S_0(B) \rightarrow T_0(B)$  iff  $f$  and  $g$  map arguments related by  $S_r(R)$  to results related by  $T_1(R)$ . Similarly, if  $V$  is the type  $\cdot \vdash \forall \alpha. S(\alpha)$ , then  $V_0$  consists of those polymorphic functions  $f$  that take a set  $A$  and return an element of  $S_0(A)$ , and also have the property that for any relation  $R \subseteq A \times B$ ,  $f(A)$  and  $f(B)$  are related by  $S_1(R)$ . Two such polymorphic functions  $f$  and  $g$  are then related by  $V_1$  iff for any relation  $R \subseteq A \times B$ ,  $f(A)$  and  $g(B)$  are related by  $S_1(R)$ . These definitions allow us to deduce interesting properties of (interpretations of) terms solely from their types. For example, for any term  $t : \forall \alpha. \alpha \rightarrow \alpha$ , the Abstraction Theorem guarantees that the interpretation  $t_0$  of  $t$  is related to itself by the relational interpretation of  $\forall \alpha. \alpha \rightarrow \alpha$ . So if we fix a set  $A$ , fix  $a \in A$ , and define a relation on  $A$  by  $R := \{(a, a)\}$ , then  $t_0(A)$  must be related to itself by the relational interpretation of  $\alpha \vdash \alpha \rightarrow \alpha$  applied to  $R$ . This means that  $t_0(A)$  must carry arguments related by  $R$  to results related by  $R$ . Since  $a$  is related to itself by  $R$ ,  $t_0(A) a$  must be related to itself by  $R$ , so that  $t_0(A) a$  must be  $a$ . That is,  $t_0$  must be the polymorphic identity function. Such applications of relational parametricity are useful in many different scenarios, e.g., when proving invariance of polymorphic functions under changes of data representation, equivalences of programs, and "free theorems" [18].

The well-known problem with Reynolds' treatment of relational parametricity (see [14]) is that the universe of sets is not impredicative, and hence the aforementioned "set"  $V_0$  cannot be formed. This issue can be resolved if we instead work in a meta-theory that has an impredicative universe; a natural choice is an extensional version of the Calculus of Inductive Constructions (CIC), i.e., a dependent type theory with a cumulative Russell-style hierarchy of universes  $\mathbb{U}_0 : \mathbb{U}_1 : \dots$ , where  $\mathbb{U}_0$  is impredicative, and extensional identity types. With this adjustment, we have two canonical relationally parametric models of System F: *i) the PER model of Longo and*

Moggi [9], internal to the theory of  $\omega$ -sets and realizable functions, and ii) Reynolds' original model<sup>1</sup>, internal to CIC.

After Reynolds' original paper, more abstract treatments of his ideas were given by, e.g., Robinson and Rosolini [15], O'Hearn and Tennent [11], Dunphy and Reddy [2], and Ghani *et al.* [5]. The approach is to use a categorical structure — reflexive graph categories for [2, 11, 15] and fibrations for [5] — to represent sets and relations, and to interpret types as appropriate functors and terms as natural transformations. In particular, [2] aims to “[address] parametricity in all its incarnations”, and similarly for [5]. Surprisingly and significantly, however, Reynolds' original model does *not* arise as a direct instance of either framework. This leads us to ask:

*What constitutes a good framework for relational parametricity?*

Our answer is that such a framework should:

1. Deliver a relationally parametric model for each instantiation of its parameters, from which it uniformly produces such models. In particular, it should allow a choice of a suitable meta-theory (the Calculus of Inductive Constructions, the theory of  $\omega$ -sets, etc.).
2. Admit the two canonical relationally parametric models mentioned above as direct instances in a natural, uniform way.
3. Abstractly formulate properties that good models of System F parametricity should be expected to satisfy.

Criterion 1 ensures that we indeed get a true *framework* rather than just a reusable blueprint for constructing models of parametricity. Criterion 2 remains unsatisfied for the frameworks of Dunphy and Reddy and of Ghani *et al.* because Reynolds' original model formulated syntactically does not satisfy certain strictness conditions imposed by [2, 5]. For example, let  $\alpha \vdash S(\alpha)$  and  $\alpha \vdash T(\alpha)$  be two types, with object interpretations  $S_0$  and  $T_0$  and relational interpretations  $S_1$  and  $T_1$ . The interpretation of the product  $\alpha \vdash S(\alpha) \times T(\alpha)$  should be an appropriate product of interpretations; that is, the object interpretation should map a set  $A$  to  $S_0(A) \times T_0(A)$  and the relational interpretation should map a relation  $R$  to  $S_1(R) \times T_1(R)$ , with the product of two relations defined in the obvious way. For the Identity Extension Lemma to hold, we need  $S_1(\text{Eq}(A)) \times T_1(\text{Eq}(A))$  to be the same as  $\text{Eq}(S_0(A) \times T_0(A))$ . Here, the equality relation  $\text{Eq}(A)$  on a set  $A$  maps  $(a, b) : A \times A$  to the type  $\text{Id}(a, b)$  of proofs of equality between  $a$  and  $b$ , so that  $a$  and  $b$  are related iff  $\text{Id}(a, b)$  is inhabited, *i.e.*, iff  $a$  is identical to  $b$ . By the induction hypothesis,  $S_1(\text{Eq}(A))$  is  $\text{Eq}(S_0(A))$ , and similarly for  $T$ , so we need to show that  $\text{Eq}(S_0(A)) \times \text{Eq}(T_0(A))$  is  $\text{Eq}(S_0(A) \times T_0(A))$ . But this is not necessarily the case since the identity type on a product is in general not *identical* to the product of identity types, but rather just suitably *isomorphic*. So the interpretation of  $\alpha \vdash S(\alpha) \times T(\alpha)$  is not necessarily an indexed or fibered functor (in the settings of [2] and [5], respectively).

Three ways to fix this problem come to mind. Firstly, we can attempt to change the meta-theory, by, e.g., imposing an additional axiom asserting that two logically equivalent propositions are definitionally equal. We do not pursue this approach here: the goal of our framework is to *directly subsume the important models in their natural meta-theories*, as per criteria 1 and 2 above, rather than require the user to augment the meta-theory with *ad hoc* axioms to

make the shoe fit. The second possibility is to use the syntactic analogue of strictification, pursued in, e.g., [1]. The idea is that instead of interpreting a closed type as a set  $A$  (on the object level), we interpret it as a set  $A$  endowed with a relation  $R_A$  that is *isomorphic*, but not necessarily identical, to the canonical discrete relation  $\text{Eq}_A$ . The chosen equality relation on the set  $A$  — more precisely, on the entire structure  $(A, (R_A, i : R_A \approx \text{Eq}_A))$  — will then be  $R_A$  rather than  $\text{Eq}_A$ . This allows us to construct  $R_A$  in a way that respects all type constructors on the nose, so that the aforementioned issue with  $\text{Eq}(S_0(A)) \times \text{Eq}(T_0(A))$  not being identical to  $\text{Eq}(S_0(A) \times T_0(A))$  is avoided. The problem, however, is that there can be many different ways to endow  $A$  with a discrete relation  $(R_A, i)$ ; in other words, the type of discrete relations on  $A$  is not contractible. It is thus unclear whether and how this “discretized” version of Reynolds' model is equivalent to the original, intended one.

Here we suggest a third approach: we record the isomorphisms witnessing the preservation of the Identity Extension Lemma for each type constructor, and propagate them through the construction. This means, however, that we can no longer interpret a type  $\alpha \vdash T(\alpha)$  as a pair of maps  $T_0 : |\text{Set}| \rightarrow \text{Set}$  and  $T_1 : |\text{Rel}| \rightarrow \text{Rel}$ ; indeed, since the domain of  $T_1$  is the discrete category  $|\text{Rel}|$ ,  $T_1$  is not required to preserve isomorphisms in  $\text{Rel}$ . As a result, even if we know that the pair  $(T_0, T_1)$  satisfies the Identity Extension Lemma, its reindexing — defined by precomposition — might not. The upshot is that the obvious “ $\lambda 2$ -fibration” corresponding to Reynolds' original model is not necessarily a fibration at all.

We solve this problem by specifying subcategories  $\mathcal{M}(0) \subseteq \text{Set}$  and  $\mathcal{M}(1) \subseteq \text{Rel}$  of *relevant isomorphisms* that form a *reflexive graph category with isomorphisms*. Abstractly, this structure gives us two *face maps* (called  $\partial_0$  and  $\partial_1$  in [2]), which represent the domain and codomain projections, and a *degeneracy* (called  $I$  in [2]), which represents the equality functor. We interpret a type  $\alpha \vdash T(\alpha)$  as a pair of functors  $T_0 : \mathcal{M}(0) \rightarrow \mathcal{M}(0)$  and  $T_1 : \mathcal{M}(1) \rightarrow \mathcal{M}(1)$  that together comprise a *face map- and degeneracy-preserving reflexive graph functor*, and interpret each term as a *face map- and degeneracy-preserving reflexive graph natural transformation*.

Since the domain of  $T_1$  is  $\mathcal{M}(1)$ ,  $T_1$  preserves all relevant isomorphisms between relations, so the reindexing of  $(T_0, T_1)$  is now well-defined. Choosing  $\mathcal{M}(1)$  to contain the isomorphism between the two relations  $\text{Eq}(S_0(A)) \times \text{Eq}(T_0(A))$  and  $\text{Eq}(S_0(A) \times T_0(A))$  yields the satisfaction of the Identity Extension Lemma for products; other type constructors follow the same pattern. We note that although the preservation of isomorphisms on the *relation* level is sufficient to carry out the model construction, we formally require the preservation of relevant isomorphisms on the *object* level, too. This makes the framework more uniform and, moreover, leads to the novel notion of a *categorical Reynolds' model*, in which interpretations of types are endowed with a functorial action on isomorphisms and all polymorphic functions respect this action. Furthermore, we go one level higher and use the ideas of Orsanigo [12] (and Ghani *et al.* [3], which it supersedes) to define a *proof-relevant categorical Reynolds' model*, in which, additionally, witnesses of relatedness are themselves suitably related via a yet higher relation.

This “2-parametric” model of course does not arise as an instance of our framework since it requires additional structure — e.g., the concept of a *2-relation* — pertaining to the higher notion of parametricity. Nevertheless, we would still like to be able to recognize it as a model parametric in the ordinary sense. Various definitions of parametricity for models of System F exist: [2, 5] are

<sup>1</sup>Since there are no set-theoretic models of System F, by the phrase “Reynolds' original model” we will always mean the version of his model that is internal to extensional CIC as described above. The need for impredicativity is inherited from Reynolds' original construction, and is not a new requirement.

examples of “internal” approaches to parametricity, where a model is considered parametric if it is produced via a specified procedure that bakes in desired features of parametricity such as the Identity Extension Lemma. On the other hand, [4, 7, 10, 15] are examples of “external” approaches to parametricity, in which reflexive graphs of models are used to endow models of interest with enough additional structure that they can reasonably be considered parametric. Surprisingly though, the proof-relevant model we give in Section 6 does not appear to satisfy any of these definitions, and in particular does not satisfy any of the external ones. The ability to construct a classifying reflexive graph seems to rely on an implicit assumption of proof-irrelevance, which we elaborate on in Section 6. However, we propose a new definition of a *relationally parametric model of System F* in Section 5 and show that it subsumes not only the two canonical parametric models of System F, but also the two novel ones we give in this paper. In particular, it subsumes the proof-relevant model given in Section 6.

The main contributions of this paper are as follows:

- We demonstrate that existing frameworks for the functorial semantics of relational parametricity for System F fail to directly subsume both canonical models of relational parametricity for System F.
- We solve this problem by developing a good abstract framework for relational parametricity that allows a choice of meta-theory, delivers both canonical relationally parametric models of System F as direct instances in a uniform way, and exposes properties that good models of System F parametricity should be expected to satisfy, e.g., guaranteeing that interpretations of terms, not just types, suitably commute with the degeneracy.
- We give a novel definition of a parametric model of System F, which is a hybrid of the external and internal approaches, and show that it subsumes both canonical models (expressed as instances of our framework).
- We give two novel relationally parametric models of System F — one of which is proof-relevant and can be seen as parametric in a higher sense (“2-parametric”) — and show that our definition recognizes both of these in a natural way, with the proof-irrelevant model arising as a direct instance of our framework.

A technical report [17] with detailed proofs is available.

## 2 Reflexive Graph Categories

Although Reynolds himself showed that his original approach to relational parametricity does not work in set theory, we can still use it as a guide for designing an abstract framework for parametricity. Instead of sets and relations, we consider abstract notions of “sets” and “relations”, and require them to be related as follows: *i*) for any relation  $R$ , there are two canonical ways of projecting an object out of  $R$ , corresponding to the domain and codomain operations, *ii*) for any object  $A$ , there is a canonical way of turning it into a relation, corresponding to the equality relation on  $A$ , and *iii*) if we start with an object  $A$ , turn it into a relation according to *ii*), and then project out an object according to *i*), we get  $A$  back. This suggests that our abstract relations and the canonical operations on them can be organized into a reflexive graph structure: categories  $\mathcal{X}_0, \mathcal{X}_1$  and functors  $f_\top, f_\perp : \mathcal{X}_1 \rightarrow \mathcal{X}_0$ ,  $\mathbf{d} : \mathcal{X}_0 \rightarrow \mathcal{X}_1$  such that  $f_\top \circ \mathbf{d} = \text{id} = f_\perp \circ \mathbf{d}$ , as is done in [2].

Since there are no set-theoretic models of System F ([14]), all of the reflexive graph structure identified above must be internal

to some ambient category  $\mathcal{C}$ . In particular,  $\mathcal{X}_0$  and  $\mathcal{X}_1$  must be categories internal in  $\mathcal{C}$ , and  $f_\top, f_\perp$ , and  $\mathbf{d}$  must be functors internal in  $\mathcal{C}$ . For Reynolds’ original model, the ambient category has types  $A : \mathbb{U}_1$  as objects and terms  $f : \Sigma_{A, B : \mathbb{U}_1} A \rightarrow B$  as morphisms. Here,  $\mathbb{U}_1$  is the universe one level above the impredicative universe  $\mathbb{U}_0$ ; we will denote  $\mathbb{U}_0$  simply by  $\mathbb{U}$  below. This ensures that  $\mathbb{U}$  is an object in  $\mathcal{C}$ . To model relations, we introduce:

$$\text{isProp}(A) := \Pi_{a, b : A} \text{Id}(a, b)$$

$$\text{Prop} := \Sigma_{A : \mathbb{U}} \text{isProp}(A)$$

The type  $\text{Prop}$  of *propositions* singles out those types in  $\mathbb{U}$  with the property that any two inhabitants, if they exist, are equal. Propositions can be used to model relations as follows: in Reynolds’ original model,  $a : A$  is related to  $b : B$  in at most one way under any relation  $R$  (either  $(a, b) \in R$  or not), so the type of proofs that  $(a, b) \in R$  is a proposition. Conversely, given  $R : A \times B \rightarrow \text{Prop}$ , we consider  $a$  and  $b$  to be related by  $R$  iff  $R(a, b)$  is inhabited.

To see the universe  $\mathbb{U}$  as a category  $\text{Set}$  internal to  $\mathcal{C}$  we take its object of objects  $\text{Set}_0$  to be  $\mathbb{U}$  and define its object of morphisms by  $\text{Set}_1 := \Sigma_{A, B : \mathbb{U}} A \rightarrow B$ . We define the category  $\mathbf{R}$  of relations by giving its objects  $\mathbf{R}_0$  and  $\mathbf{R}_1$  of objects and morphisms, respectively:

$$\mathbf{R}_0 := \Sigma_{A, B : \text{Set}} A \times B \rightarrow \text{Prop}$$

$$\mathbf{R}_1 := \Sigma_{((A_1, A_2), R_A), ((B_1, B_2), R_B) : \mathbf{R}_0} \Sigma_{(f, g) : (A_1 \rightarrow B_1) \times (A_2 \rightarrow B_2)} \Pi_{(a_1, a_2) : A_1 \times A_2} R_A(a_1, a_2) \rightarrow R_B(f(a_1), g(a_2))$$

We clearly have two internal functors from  $\mathbf{R}$  to  $\text{Set}$  corresponding to the domain and codomain projections, respectively. We also have an internal functor  $\text{Eq}$  from  $\text{Set}$  to  $\mathbf{R}$  that constructs an equality relation with  $\text{Eq } A := ((A, A), \text{Id}_A)$  and  $\text{Eq } ((A, B), f) := ((\text{Eq } A, \text{Eq } B), (f, f), \text{ap}_f)$ . Here, the term  $\text{ap}_f$  of type  $\text{Id}_A(a_1, a_2) \rightarrow \text{Id}_B(f(a_1), f(a_2))$  is defined as usual by  $\text{Id}$ -induction and witnesses the fact that  $f$  respects equality.

These observations motivate the next two definitions, in which we denote the category of categories and functors internal to  $\mathcal{C}$  by  $\text{Cat}(\mathcal{C})$ , and assume  $\mathcal{C}$  is locally small and has all finite products. (A category is *locally small* if each of its hom-sets is small, i.e., is a set rather than a proper class.)

**Definition 2.1.** A reflexive graph structure  $\mathcal{X}$  on a category  $\mathcal{C}$  consists of:

- objects  $\mathcal{X}(0)$  and  $\mathcal{X}(1)$  of  $\mathcal{C}$
- distinct arrows  $\mathcal{X}(f_\star) : \mathcal{X}(1) \rightarrow \mathcal{X}(0)$  for  $\star : \text{Bool}$
- an arrow  $\mathcal{X}(\mathbf{d}) : \mathcal{X}(0) \rightarrow \mathcal{X}(1)$

such that  $\mathcal{X}(f_\star) \circ \mathcal{X}(\mathbf{d}) = \text{id}$ .

The requirement that the two face maps  $\mathcal{X}(f_\top)$  and  $\mathcal{X}(f_\perp)$  are distinct is to ensure that there are enough relations for the notion of relation-preservation to be meaningful. Otherwise, as also observed in [2], we could see *any* category  $\mathcal{C}$  as supporting a trivial reflexive graph structure whose only relations are the equality ones. For readers familiar with [7], the condition  $\mathcal{X}(f_\top) \neq \mathcal{X}(f_\perp)$  serves a purpose similar to that of the requirement in Definition 8.6.2 of [7] that the fiber category  $\mathbb{F}_1$  over the terminal object in  $\mathbb{C}$  is the category of relations in the preorder fibration  $\mathbb{D} \rightarrow \mathbb{E}$  on the fiber category  $\mathbb{E}_1$  over the terminal object in  $\mathbb{B}$ . Both conditions imply that some relations must be *heterogeneous*. But while in [7] relations are obtained in a standard way as predicates (given by a preorder fibration) over a product, we do not assume that relations are constructed in any

specific way, but rather only that the abstract operations on relations suitably interact. Moreover, since the two face maps  $\mathcal{X}(f_\top)$  and  $\mathcal{X}(f_\perp)$  are distinct, any morphism generated by the face maps and the degeneracy  $\mathcal{X}(\mathbf{d})$  must be one of the seven distinct maps  $\text{id}_{\mathcal{X}(0)}$ ,  $\text{id}_{\mathcal{X}(1)}$ ,  $\mathcal{X}(f_\star)$ ,  $\mathcal{X}(\mathbf{d})$ , and  $\mathcal{X}(\mathbf{d}) \circ \mathcal{X}(f_\star)$  for  $\star : \text{Bool}$ . Every such morphism thus has a canonical representation.

**Definition 2.2.** A reflexive graph category (on  $C$ ) is a reflexive graph structure on  $\text{Cat}(C)$ .

**Example 2.3** (PER model). We take the ambient category  $C$  to be the category of  $\omega$ -sets, given in Definition 6.3 of [9]. We construct a reflexive graph category, which we call  $\mathcal{R}_{\text{PER}}$ , as follows. The internal category  $\mathcal{R}_{\text{PER}}(0)$  of “sets” is the category  $\mathbf{M}'$  given in Definition 8.4 of [9]. Informally, the objects of  $\mathbf{M}'$  are partial equivalence relations on  $\mathbb{N}$ , and the morphisms are realizable functions that respect such relations. To define the internal category  $\mathcal{R}_{\text{PER}}(1)$  of “relations”, we first construct its object of objects. The carrier of this  $\omega$ -set is the set of pairs of the form  $R := ((A_d, A_c), R_A)$ , where  $A_d$  and  $A_c$  are partial equivalence relations and  $R_A$  is a saturated predicate on the product  $\text{PER } A_d \times A_c$ . A saturated predicate on a PER  $A$  is a predicate on  $\mathbb{N}$  such that  $a_1 \sim_A a_2$  and  $R(a_1)$  imply  $R(a_2)$ . To finish the construction of our object of objects for  $\mathcal{R}_{\text{PER}}(1)$  we take any pair  $((A_d, A_c), R_A)$  as above to be realized by any natural number.

The carrier of the object of morphisms for  $\mathcal{R}_{\text{PER}}(1)$  comprises all pairs of the form

$$(((A_d, A_c), R_A), ((B_d, B_c), R_B)), (\{m_1\}_{A_d \rightarrow B_d}, \{m_2\}_{A_c \rightarrow B_c}))$$

satisfying the condition that, for any  $k, l$  such that  $k \sim_{A_d} l$ ,  $l \sim_{A_c} l$ , and  $R_A((k, l))$  holds,  $R_B((m_1 \cdot k, m_2 \cdot l))$  holds as well. The first component records the domain and codomain of the morphism and the second component is a pair of equivalence classes under the specified exponential PERs. As in [9], we denote the application of the  $n^{\text{th}}$  partial recursive function to a natural number  $a$  in its domain by  $n \cdot a$ . To finish the construction of the object of morphisms for  $\mathcal{R}_{\text{PER}}(1)$ , we take a pair of pairs as above to be realized by a natural number  $k$  iff  $\text{fst}(k) \sim_{A_d \rightarrow B_d} m_1$  and  $\text{snd}(k) \sim_{A_c \rightarrow B_c} m_2$ .

We again have two internal functors  $\mathcal{R}_{\text{PER}}(f_\top)$  and  $\mathcal{R}_{\text{PER}}(f_\perp)$  from  $\mathcal{R}_{\text{PER}}(1)$  to  $\mathcal{R}_{\text{PER}}(0)$  corresponding to the two projections. We also have an equality functor  $\text{Eq}$  from  $\mathcal{R}_{\text{PER}}(0)$  to  $\mathcal{R}_{\text{PER}}(1)$  whose action on objects is given by  $\text{Eq } A := ((A, A), \Delta_A)$ , where  $\Delta_A(k)$  iff  $\text{fst}(k) \sim_A \text{snd}(k)$ , and whose action on morphisms is given by

$$\text{Eq}((A, B), \{m\}_{A \rightarrow B}) := ((\text{Eq } A, \text{Eq } B), (\{m\}_{A \rightarrow B}, \{m\}_{A \rightarrow B}))$$

**Example 2.4** (Reynolds’ model). We obtain a reflexive graph category  $\mathcal{R}_{\text{REY}}$  by taking  $\mathcal{R}_{\text{REY}}(0) := \text{Set}$ ,  $\mathcal{R}_{\text{REY}}(1) := \text{R}$ , and  $\mathcal{R}_{\text{REY}}(\mathbf{d}) := \text{Eq}$ , and letting  $\mathcal{R}_{\text{REY}}(f_\top)$  and  $\mathcal{R}_{\text{REY}}(f_\perp)$  be the functors corresponding to the domain and codomain projections, respectively.

If  $\mathcal{X}$  is a reflexive graph category, then the discrete graph category  $|\mathcal{X}|$  and the product reflexive graph category  $\mathcal{X}^n$  for  $n \in \mathbb{N}$  are defined in the obvious ways:  $|\mathcal{X}(l)|$  has the same objects as  $\mathcal{X}(l)$  but only the identity morphisms, and  $(\mathcal{X} \times \mathcal{Y})(l) = \mathcal{X}(l) \times \mathcal{Y}(l)$  for  $l \in \{0, 1\}$ . For the latter, the product on the right-hand side is a product of internal categories, which exists because  $C$  has finite products by assumption. To simplify the presentation, we will omit explicit mentions of the category  $C$ , and treat definitions and constructions internal to  $C$  as though they were external.

Given a reflexive graph category  $\mathcal{X}$  axiomatizing the sets and relations, an obvious first attempt at pushing Reynolds’ original

idea through is to take the interpretation  $\llbracket T \rrbracket$  of a type  $\bar{a} \vdash T$  with  $n$  free type variables to be a pair  $(\llbracket T \rrbracket(0), \llbracket T \rrbracket(1))$ , where  $\llbracket T \rrbracket(0) : |\mathcal{X}(0)|^n \rightarrow \mathcal{X}(0)$  and  $\llbracket T \rrbracket(1) : |\mathcal{X}(1)|^n \rightarrow \mathcal{X}(1)$  are functions giving the “set” and “relation” interpretations of the type  $T$ . Although as explained in the introduction, this approach will need some tweaking – we will need to endow  $\llbracket T \rrbracket(0)$  and  $\llbracket T \rrbracket(1)$  with actions on some morphisms – it suggests:

**Definition 2.5.** Let  $\mathcal{X}$  and  $\mathcal{Y}$  be reflexive graph categories. A reflexive graph functor  $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$  is a pair  $(\mathcal{F}(0), \mathcal{F}(1))$  of functors such that  $\mathcal{F}(0) : \mathcal{X}(0) \rightarrow \mathcal{Y}(0)$  and  $\mathcal{F}(1) : \mathcal{X}(1) \rightarrow \mathcal{Y}(1)$ .

Writing  $T_0$  for  $\llbracket T \rrbracket(0)$  and  $T_1$  for  $\llbracket T \rrbracket(1)$ , we recall from the introduction that  $T_0$  and  $T_1$  should be appropriately related via the domain and codomain projections and the equality functor. Since the two face maps  $\mathcal{X}(f_\star)$  now model the projections, and the degeneracy  $\mathcal{X}(\mathbf{d})$  models the equality functor, we end up with the following conditions: *i*) for each object  $\bar{R}$  in  $\mathcal{X}(1)^n$ , we have  $\mathcal{X}(f_\star) T_1(\bar{R}) = T_0(\mathcal{X}(f_\star)^n \bar{R})$ , and *ii*) for each object  $\bar{A}$  in  $\mathcal{X}(0)^n$ , we have  $\mathcal{X}(\mathbf{d}) T_0(\bar{A}) = T_1(\mathcal{X}(\mathbf{d})^n \bar{A})$ . We examine what these conditions imply for Reynolds’ model by considering the product  $\alpha \vdash S(\alpha) \times T(\alpha)$  of two types  $\alpha \vdash S(\alpha)$  and  $\alpha \vdash T(\alpha)$ . By the induction hypothesis,  $S$  and  $T$  are interpreted as pairs  $(S_0, S_1)$  and  $(T_0, T_1)$ , where  $S_0, T_0 : \text{Set}_0 \rightarrow \text{Set}_0$  and  $S_1, T_1 : \text{R}_0 \rightarrow \text{R}_0$  satisfy *i*) and *ii*). The interpretation of a product should be a product of interpretations, i.e.,  $(S \times T)_0 A := S_0(A) \times T_0(A)$  and  $(S \times T)_1 R := S_1(R) \times T_1(R)$ . It remains to be seen that this interpretation satisfies *i*) and *ii*).

Fix a relation  $R$  on  $A$  and  $B$ . Condition *i*) entails that  $S_1(R) := ((S_0(A), S_0(B)), R_S)$  and  $T_1(R) := ((T_0(A), T_0(B)), R_T)$  for some  $R_S$  and  $R_T$ . Thus  $S_1(R) \times T_1(R)$  has the form  $((S_0(A) \times T_0(A), S_0(B) \times T_0(B)), R_{S \times T})$ , where  $R_{S \times T}$  maps a pair of pairs  $((a, b), (c, d))$  to  $R_S(a, c) \times R_T(b, d)$ . Thus *i*) is satisfied simply by construction, which leads us to define:

**Definition 2.6.** A reflexive graph functor  $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$  is face map-preserving if  $\mathcal{Y}(f_\star) \circ \mathcal{F}(1) = \mathcal{F}(0) \circ \mathcal{X}(f_\star)$  for every  $\star \in \text{Bool}$ .

In Reynolds’ model, condition *ii*) gives that  $S_1(\text{Eq}(A))$  is  $\text{Eq}(S_0(A))$  for any set  $A$ , and similarly for  $T$ . We thus need to show that  $\text{Eq}(S_0(A)) \times \text{Eq}(T_0(A))$  is  $\text{Eq}(S_0(A) \times T_0(A))$ . But while the domains and codomains of these two relations agree (all are  $S_0(A) \times T_0(A)$ ), the former maps  $((a, b), (c, d))$  to  $\text{Id}(a, c) \times \text{Id}(b, d)$ , while the latter maps it to  $\text{Id}((a, b), (c, d))$ . These two types are not necessarily identical, but they are *isomorphic* (i.e., there are functions back and forth that compose to identity on both sides).

We thus relax condition *ii*) to allow an isomorphism  $\varepsilon_T(\bar{A}) : \mathcal{X}(\mathbf{d}) T_0(\bar{A}) \cong T_1(\mathcal{X}(\mathbf{d})^n \bar{A})$ . In fact, we can require more: since the domains and codomains of  $\mathcal{X}(\mathbf{d}) T_0(\bar{A})$  and  $T_1(\mathcal{X}(\mathbf{d})^n \bar{A})$  coincide by condition *i*), we can insist that both projections map the isomorphism  $\varepsilon_T(\bar{A})$  to the identity morphism on  $T_0(\bar{A})$ . This coherence condition is a natural counterpart to the equation  $\mathcal{X}(f_\star) \circ \mathcal{X}(\mathbf{d}) = \text{id}$ , and turns out to be not just a design choice but a necessary requirement: in Reynolds’ model, for instance, the proof that the interpretations of  $\forall$ -types (as defined later) suitably commute with the functor  $\text{Eq}$  depends precisely on the morphisms underlying the maps  $\varepsilon_T(\bar{A})$  being identities. This suggests:

**Definition 2.7.** A reflexive graph functor  $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$  is said to be degeneracy-preserving if there is a natural isomorphism  $\varepsilon_{\mathcal{F}} : \mathcal{Y}(\mathbf{d}) \circ \mathcal{F}(0) \rightarrow \mathcal{F}(1) \circ \mathcal{X}(\mathbf{d})$  satisfying the coherence condition  $\mathcal{Y}(f_\star) \circ \varepsilon_{\mathcal{F}} = \text{id}$ .

As a first approximation, we can try to interpret a type  $\bar{\alpha} \vdash T$  with  $n$  free type variables as a face map- and degeneracy-preserving reflexive graph functor  $(T_0, T_1) : |\mathcal{X}|^n \rightarrow \mathcal{X}$ . Reynolds' original idea for interpreting terms suggests that the interpretation of a term  $\bar{\alpha}; x : S \vdash t : T$  should be a (vacuously) natural transformation  $t_0 : S_0 \rightarrow T_0$ . As observed in [5], the Abstraction Theorem can then be formulated as follows: there is a (vacuously) natural transformation  $t_1 : S_1 \rightarrow T_1$  such that, for any object  $\bar{R}$  in  $\mathcal{X}(1)^n$ , we have  $\mathcal{X}(f_\star) t_1(\bar{R}) = t_0(\mathcal{X}(f_\star)^n \bar{R})$ . To see that this does indeed give what we want, we revisit Reynolds' model. There, the face maps are the domain and codomain projections and an object  $\bar{R}$  in  $\mathcal{X}(1)^n$  is an  $n$ -tuple of relations. Denote  $\mathcal{X}(f_\top)^n \bar{R}$  by  $\bar{A}$  and  $\mathcal{X}(f_\perp)^n \bar{R}$  by  $\bar{B}$ . Then  $t_1(\bar{R})$  is a morphism of relations from  $S_1(\bar{R})$  to  $T_1(\bar{R})$  and, since  $S_1$  and  $T_1$  are face map-preserving,  $S_1(\bar{R}) := ((S_0(\bar{A}), S_0(\bar{B})), R_S)$  and  $T_1(\bar{R}) := ((T_0(\bar{A}), T_0(\bar{B})), R_T)$  for some  $R_S$  and  $R_T$ . By definition,  $t_1(\bar{R})$  gives maps  $f : S_0(\bar{A}) \rightarrow T_0(\bar{A})$ ,  $g : S_0(\bar{B}) \rightarrow T_0(\bar{B})$ , together with a map  $h : \prod_{(a_1, a_2) : S_0(\bar{A}) \times S_0(\bar{B})} R_S(a_1, a_2) \rightarrow R_T(f(a_1), g(a_2))$  stating precisely that  $f$  and  $g$  map related inputs to related outputs. By definition,  $\mathcal{X}(f_\top) t_1(\bar{R})$  is  $((S_0(\bar{A}), T_0(\bar{A})), f)$  and  $\mathcal{X}(f_\perp) t_1(\bar{R})$  is  $((S_0(\bar{B}), T_0(\bar{B})), g)$ , so the condition that  $\mathcal{X}(f_\star) t_1(\bar{R})$  is  $t_0(\mathcal{X}(f_\star)^n \bar{R})$  implies that the maps underlying  $t_0(\bar{A})$  and  $t_0(\bar{B})$  must be  $f$  and  $g$ , respectively, and so must indeed map related inputs to related outputs, as witnessed by  $h$ . Pairing the natural transformations  $t_0$  and  $t_1$  motivates:

**Definition 2.8.** Let  $\mathcal{F}, \mathcal{G} : \mathcal{X} \rightarrow \mathcal{Y}$  be reflexive graph functors. A reflexive graph natural transformation  $\eta : \mathcal{F} \rightarrow \mathcal{G}$  is a pair  $(\eta(0), \eta(1))$  of natural transformations  $\eta(0) : \mathcal{F}(0) \rightarrow \mathcal{G}(0)$  and  $\eta(1) : \mathcal{F}(1) \rightarrow \mathcal{G}(1)$ .

The Abstraction Theorem then further suggests defining:

**Definition 2.9.** A reflexive graph natural transformation  $\eta : \mathcal{F} \rightarrow \mathcal{G}$  is *face map-preserving* if  $\mathcal{F}$  and  $\mathcal{G}$  are face map-preserving and, for each  $\star \in \mathbf{Bool}$ , we have  $\mathcal{Y}(f_\star) \circ \eta(1) = \eta(0) \circ \mathcal{X}(f_\star)$ .

The interpretation of a term  $\bar{\alpha}; x : S \vdash t : T$  should then be a face map-preserving natural transformation from  $(S_0, S_1)$  to  $(T_0, T_1)$ . We also have the dual notion:

**Definition 2.10.** A reflexive graph natural transformation  $\eta : \mathcal{F} \rightarrow \mathcal{G}$  is *degeneracy-preserving* if  $\mathcal{F}$  and  $\mathcal{G}$  are degeneracy-preserving, as witnessed by the natural isomorphisms  $\varepsilon_{\mathcal{F}}$  and  $\varepsilon_{\mathcal{G}}$ , respectively, and, for every  $X$  in  $\mathcal{X}(0)$ , we have  $(\eta(1) (\mathcal{X}(\mathbf{d}) X)) \circ \varepsilon_{\mathcal{F}}(X) = \varepsilon_{\mathcal{G}}(X) \circ (\mathcal{Y}(\mathbf{d}) (\eta(0) X))$ .

There is no explicit analogue of Definition 2.10 in Reynolds' model for the following reason: Reynolds' model (as well as the PER model) is proof-irrelevant, in the precise sense that the functor  $\langle \mathcal{X}(f_\perp), \mathcal{X}(f_\top) \rangle$  is faithful, and this condition is sufficient to guarantee that *any* face map-preserving natural transformation is automatically degeneracy-preserving as well. This may or may not be the case in proof-relevant models (although in the model from Section 6 it is), so we explicitly restrict attention below only to those natural transformations that are face map- and degeneracy-preserving (as also done in [2]), and omit further mention of these properties.

Identity and composition for natural transformations between reflexive graph functors are defined levelwise. Identity for reflexive graph functors is also obvious, but composition requires some care:

**Definition 2.11.** For reflexive graph functors  $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$  and  $\mathcal{G} : \mathcal{Y} \rightarrow \mathcal{Z}$ , the reflexive graph functor  $\mathcal{G} \circ \mathcal{F} : \mathcal{X} \rightarrow \mathcal{Z}$  is defined as follows:

- $(\mathcal{G} \circ \mathcal{F})(l) := \mathcal{G}(l) \circ \mathcal{F}(l)$
- $\varepsilon_{\mathcal{G} \circ \mathcal{F}}(X) := (\mathcal{G}(1) \varepsilon_{\mathcal{F}}) \circ \varepsilon_{\mathcal{G}}(\mathcal{F}(0) X)$

Here, the first composition is a composition of functors and the second is a composition of morphisms in the category  $\mathcal{Z}(1)$ .

Given reflexive graph functors  $\mathcal{F}_1, \mathcal{F}_2 : \mathcal{X} \rightarrow \mathcal{Y}$  and  $\mathcal{G}_1, \mathcal{G}_2 : \mathcal{Y} \rightarrow \mathcal{Z}$ , and natural transformations  $\varepsilon : \mathcal{F}_1 \rightarrow \mathcal{F}_2$  and  $\eta : \mathcal{G}_1 \rightarrow \mathcal{G}_2$ , the compositions  $\eta \circ \mathcal{F}_1 : \mathcal{G}_1 \circ \mathcal{F}_1 \rightarrow \mathcal{G}_2 \circ \mathcal{F}_1$  and  $\mathcal{G}_1 \circ \varepsilon : \mathcal{G}_1 \circ \mathcal{F}_1 \rightarrow \mathcal{G}_1 \circ \mathcal{F}_2$  are defined levelwise in the obvious way. A particular reflexive graph functor of interest, which we will use to interpret type variables, is projection:

**Definition 2.12.** Given a reflexive graph category  $\mathcal{X}$  and  $i \in \{1, \dots, n\}$ , the  $i^{\text{th}}$  reflexive graph projection functor is the reflexive graph functor  $\text{pr}_i^n : \mathcal{X}^n \rightarrow \mathcal{X}$ , where  $\text{pr}_i^n(l) : \mathcal{X}(l)^n \rightarrow \mathcal{X}(l)$  is the usual  $i^{\text{th}}$  projection functor and  $\varepsilon_{\text{pr}_i^n}(X) := \text{id}$ .

Dually, we have the following:

**Definition 2.13.** For reflexive graph functors  $\mathcal{F}_1, \dots, \mathcal{F}_m : \mathcal{X} \rightarrow \mathcal{Y}$ , the reflexive graph functor  $\langle \mathcal{F}_1, \dots, \mathcal{F}_m \rangle : \mathcal{X} \rightarrow \mathcal{Y}^m$  is defined as follows:

- $\langle \mathcal{F}_1, \dots, \mathcal{F}_m \rangle(l) := \langle \mathcal{F}_1(l), \dots, \mathcal{F}_m(l) \rangle$
- $\varepsilon_{\langle \mathcal{F}_1, \dots, \mathcal{F}_m \rangle}(X) := \langle \varepsilon_{\mathcal{F}_1}(X), \dots, \varepsilon_{\mathcal{F}_m}(X) \rangle$

Similarly, given reflexive graph natural transformations  $\eta_1 : \mathcal{F}_1 \rightarrow \mathcal{G}_1, \dots, \eta_m : \mathcal{F}_m \rightarrow \mathcal{G}_m$ , the reflexive graph natural transformation  $\langle \eta_1, \dots, \eta_m \rangle : \langle \mathcal{F}_1, \dots, \mathcal{F}_m \rangle \rightarrow \langle \mathcal{G}_1, \dots, \mathcal{G}_m \rangle$  is defined in the obvious way.

### 3 Reflexive Graph Categories with Isomorphisms

As noted above, if we try to interpret a type  $\bar{\alpha} \vdash T$  as a reflexive graph functor  $\llbracket T \rrbracket : \mathcal{X}^n \rightarrow \mathcal{X}$  we encounter a problem with contravariance. Specifically, if  $\alpha \vdash A$  and  $\alpha \vdash B$  are types, then to interpret the function type  $\alpha \vdash A \rightarrow B$  as the exponential of  $\llbracket A \rrbracket$  and  $\llbracket B \rrbracket$ ,  $\llbracket A \rightarrow B \rrbracket(0)$  must map each object  $X$  to the exponential  $(\llbracket A \rrbracket(0) X) \Rightarrow (\llbracket B \rrbracket(0) X)$  and each morphism  $f : X \rightarrow Y$  to a morphism from  $(\llbracket A \rrbracket(0) X) \Rightarrow (\llbracket B \rrbracket(0) X)$  to  $(\llbracket A \rrbracket(0) Y) \Rightarrow (\llbracket B \rrbracket(0) Y)$ . But there is no canonical way to construct a morphism of this type because  $\llbracket A \rrbracket(0) f$  goes in the wrong direction. This is a well-known problem that is unrelated to parametricity.

The usual solution is to require the domains of the functors interpreting types to be discrete, so that  $\llbracket T \rrbracket : |\mathcal{X}|^n \rightarrow \mathcal{X}$ . However, as noted in the introduction, this will not work in our setting. Consider types  $\alpha \vdash S(\alpha)$  and  $\cdot \vdash T$ . By the induction hypothesis,  $\llbracket S \rrbracket : |\mathcal{X}| \rightarrow \mathcal{X}$  and  $\llbracket T \rrbracket : 1 \rightarrow \mathcal{X}$  are face map- and degeneracy-preserving reflexive graph functors. The interpretation of the type  $\cdot \vdash S[\alpha := T]$  should be given by the composition  $\llbracket S \rrbracket \circ \llbracket T \rrbracket : 1 \rightarrow \mathcal{X}$ , which should be a face map- and degeneracy-preserving functor. While preservation of face maps is easy to prove, preservation of degeneracies poses a problem: writing  $S_0$  and  $S_1$  for  $\llbracket S \rrbracket(0)$  and  $\llbracket S \rrbracket(1)$ , and similarly for  $T$ , we need  $S_1(T_1)$  to be isomorphic to the degeneracy  $\mathbf{d}(S_0(T_0))$ . By assumption,  $T_1$  is isomorphic to the degeneracy  $\mathbf{d}(T_0)$ , and  $S_1(\mathbf{d}(T_0))$  is isomorphic to  $\mathbf{d}(S_0(T_0))$ , so if we knew that  $S_1$  mapped isomorphic relations to isomorphic relations we would be done. But since the domain of  $S_1$  is  $|\mathcal{X}(1)|$ , there is no reason that it should preserve non-identity isomorphisms of  $\mathcal{X}(1)$ .

In this paper we solve this contravariance problem in a different way. We first note that the issue does not arise if  $\llbracket A \rrbracket(0) f$  is an isomorphism, even if that isomorphism is not the identity. This leads us to require, for each  $l \in \{0, 1\}$ , a wide subcategory  $\mathcal{M}(l) \subseteq \mathcal{X}(l)$  such that every morphism in  $\mathcal{M}(l)$  is in fact an isomorphism.

**Definition 3.1.** Given a reflexive graph category  $\mathcal{X}$ , a *reflexive graph subcategory* of  $\mathcal{X}$  is a reflexive graph category  $\mathcal{M}$  together with a reflexive graph functor  $\mathcal{I} : \mathcal{M} \rightarrow \mathcal{X}$  such that:

- The object and morphism parts of  $\mathcal{I}$  are monomorphisms.
- $\mathcal{I}(0) \circ \mathcal{M}(f_\star) = \mathcal{X}(f_\star) \circ \mathcal{I}(1)$  for  $\star \in \mathbf{Bool}$ .
- $\mathcal{I}(1) \circ \mathcal{M}(d) = \mathcal{X}(d) \circ \mathcal{I}(1)$ .

The subcategory  $(\mathcal{M}, \mathcal{I})$  is *wide* if the object parts of  $\mathcal{I}(0)$  and  $\mathcal{I}(1)$  are isomorphisms.

The last two conditions in Definition 3.1 guarantee that  $\mathcal{I}$  preserves face maps and degeneracies on the nose.

**Definition 3.2.** A *reflexive graph category with isomorphisms* is a reflexive graph category  $\mathcal{X}$  together with a wide reflexive graph subcategory  $(\mathcal{M}, \mathcal{I})$  such that every morphism in  $\mathcal{M}(l)$ ,  $l \in \{0, 1\}$ , is an isomorphism.

We view  $\mathcal{M}(l)$  as selecting the *relevant isomorphisms* of  $\mathcal{X}(l)$ , in the sense that a morphism of  $\mathcal{X}(l)$  is relevant iff it lies in the image of  $\mathcal{I}(l)$ . Given a reflexive graph category with isomorphisms  $(\mathcal{X}, (\mathcal{M}, \mathcal{I}))$  we can now interpret a type  $\bar{\alpha} \vdash T$  with  $n$  free type variables as a reflexive graph functor  $\llbracket T \rrbracket : \mathcal{M}^n \rightarrow \mathcal{M}$ . It is important that  $\llbracket T \rrbracket$  carries (tuples of) relevant isomorphisms to relevant isomorphisms: if  $\llbracket T \rrbracket$  were instead a functor from  $\mathcal{M}^n$  to  $\mathcal{X}$ , then it would not be possible to define substitution (see Definition 4.2).

A trivial choice is to take  $\mathcal{M} := |\mathcal{X}|$ . Then  $\llbracket T \rrbracket : |\mathcal{X}|^n \rightarrow |\mathcal{X}|$  and  $\varepsilon_{\llbracket T \rrbracket}$  is necessarily the identity natural transformation, so  $\llbracket T \rrbracket$  preserves degeneracies on the nose. This instantiation shows that, despite being motivated by Reynolds' model, for which the Identity Extension Lemma holds only up to isomorphism, our framework can also uniformly subsume strict models of parametricity, for which the Identity Extension Lemma holds on the nose.

**Example 3.3** (PER model, continued). We take  $\mathcal{M} := |\mathcal{R}_{\text{PER}}|$ .

**Example 3.4** (A categorical version of Reynolds' model). For each  $l$ , we take the objects of  $\mathcal{M}(l)$  to be the objects of  $\mathcal{R}_{\text{REY}}(l)$ , and the morphisms of  $\mathcal{M}(l)$  to be *all* isomorphisms of  $\mathcal{R}_{\text{REY}}(l)$ . For example, the morphisms of  $\mathcal{M}(0)$  are

$$\{(i, j) : \text{Set}_1 \times \text{Set}_1 \ \& \ i_d = j_c \times i_c = j_d \times j \circ i = \text{id} \times i \circ j = \text{id}\}$$

Here and at several places below we write  $a = b$  for  $\text{Id}(a, b)$  and  $\{x : A \ \& \ B(x)\}$  for  $\Sigma_{x:A} B(x)$  to enhance readability. Moreover,  $\circ$  and  $\text{id}$  are composition and identity in the category  $\text{Set}$ , and we use the subscripts  $(\cdot)_d$  and  $(\cdot)_c$  to denote the domain and codomain of a morphism. The first (or second) projection gives the required mono from  $\mathcal{M}(0)$  to  $\text{Set}_1$ . We denote the resulting reflexive graph category with isomorphisms by  $\mathcal{R}_{\text{CREY}}$ .

**Example 3.5** (Reynolds' model, continued). As mentioned in the introduction, to push the constructions through it is sufficient to require preservation of isomorphisms on the relation level only. This means that on the set level, we can take the relevant isomorphisms to be just the identities, *i.e.*,  $\mathcal{M}(0) := |\mathcal{R}_{\text{REY}}(0)|$ . On the relation level, we take the objects of  $\mathcal{M}(1)$  to be the objects of  $\mathcal{R}_{\text{REY}}(1)$  — *i.e.*, we take all relations — and the morphisms of  $\mathcal{M}(1)$  to be those

isomorphisms of  $\mathcal{R}_{\text{REY}}(1)$  whose images under the two face maps are identities (this last condition is necessary since face maps must preserve relevant isomorphisms). Specifically, the morphisms of  $\mathcal{M}(1)$  are

$$\{(i, j) : R_1 \times R_1 \ \& \ i_d = j_c \times i_c = j_d \times j \circ i = \text{id} \times i \circ j = \text{id} \times i_\perp = \text{id}\}$$

Here, we use the subscripts  $(\cdot)_\top$  and  $(\cdot)_\perp$  to denote the image of a morphism in  $R_1$  under the corresponding face map.

With this infrastructure in place we can now interpret a term  $\bar{\alpha}; x : S \vdash t : T$  as a natural transformation from  $\mathcal{I} \circ \llbracket S \rrbracket$  to  $\mathcal{I} \circ \llbracket T \rrbracket$ . Importantly, the components of such a natural transformation are drawn from  $\mathcal{X}(l)$  (as witnessed by post-composition with  $\mathcal{I}$ ), rather than just  $\mathcal{M}(l)$ , as would be the case if we interpreted  $t$  as a natural transformation from  $\llbracket S \rrbracket$  to  $\llbracket T \rrbracket$ . In fact, this latter interpretation would not even be sensible, since not every term gives rise to an isomorphism (most do not).

## 4 Cartesian Closed Reflexive Graph Categories with Isomorphisms

We want to interpret a type context of length  $n$  as the natural number  $n$ , types with  $n$  free type variables as reflexive graph functors from  $\mathcal{M}^n$  to  $\mathcal{M}$ , and terms with  $n$  free type variables as natural transformations between reflexive graph functors with codomain  $\mathcal{X}$ . Following the standard procedure, we first define, for each  $n$ , a category  $\mathcal{M}^n \rightarrow \mathcal{M}$  to interpret expressions with  $n$  free type variables, and then combine these categories using the usual Grothendieck construction. This gives a fibration whose fiber over  $n$  is  $\mathcal{M}^n \rightarrow \mathcal{M}$ .

**Definition 4.1.** The category  $\mathcal{M}^n \rightarrow \mathcal{M}$  has as its objects the face map- and degeneracy-preserving reflexive graph functors from  $\mathcal{M}^n$  to  $\mathcal{M}$ , and as its morphisms from  $\mathcal{F}$  to  $\mathcal{G}$  the face map- and degeneracy-preserving reflexive graph natural transformations from  $\mathcal{I} \circ \mathcal{F}$  to  $\mathcal{I} \circ \mathcal{G}$ .

If  $\mathcal{F}$  and  $\mathcal{G}$  are degeneracy-preserving then  $\mathcal{I} \circ \mathcal{F}$  and  $\mathcal{I} \circ \mathcal{G}$  are as well, and it is therefore sensible to require natural transformations between the latter two to be degeneracy-preserving. To move between the fibers we need a notion of substitution:

**Definition 4.2.** For any  $m$ -tuple  $\mathbf{F} := (F_1, \dots, F_m)$  of objects in  $\mathcal{M}^n \rightarrow \mathcal{M}$ , the functor  $\mathbf{F}^*$  from  $\mathcal{M}^m \rightarrow \mathcal{M}$  to  $\mathcal{M}^n \rightarrow \mathcal{M}$  is defined by  $\mathbf{F}^*(G) := G \circ \langle F_1, \dots, F_m \rangle$  for objects and  $\mathbf{F}^*(\eta) := \eta \circ \langle F_1, \dots, F_m \rangle$  for morphisms.

When giving a categorical interpretation of System F, a category for interpreting type contexts is also required. Writing  $\mathcal{R}$  for the tuple  $(\mathcal{X}, (\mathcal{M}, \mathcal{I}))$ , we define:

**Definition 4.3.** The *category of contexts*  $\text{Ctx}(\mathcal{R})$  is given by:

- objects are natural numbers
- morphisms from  $n$  to  $m$  are  $m$ -tuples of objects in  $\mathcal{M}^n \rightarrow \mathcal{M}$
- the identity  $\text{id}_n : n \rightarrow n$  has as its  $i^{\text{th}}$  component the  $i^{\text{th}}$  projection functor  $\text{pr}_i^n$
- given morphisms  $\mathbf{F} : n \rightarrow m$  and  $\mathbf{G} = (G_1, \dots, G_k) : m \rightarrow k$ , the  $i^{\text{th}}$  component of the composition  $\mathbf{G} \circ \mathbf{F} : n \rightarrow k$  is  $\mathbf{F}^*(G_i)$

Defining the product  $n \times 1$  in  $\text{Ctx}(\mathcal{R})$  to be the natural number sum  $n + 1$  shows that  $\text{Ctx}(\mathcal{R})$  can model System F type contexts:

**Lemma 4.4.** *The category  $\text{Ctx}(\mathcal{R})$  has a terminal object 0 and products  $(-) \times 1$ .*

The categories  $\text{Ctx}(\mathcal{R})$  and  $\mathcal{M}^n \rightarrow \mathcal{M}$  can be combined to give:

**Definition 4.5.** The category  $\int_n \mathcal{M}^n \rightarrow \mathcal{M}$  is defined as follows:

- objects are pairs  $(n, F)$ , where  $F$  is an object in  $\mathcal{M}^n \rightarrow \mathcal{M}$
- morphisms from  $(n, F)$  to  $(m, G)$  are pairs  $(F, \eta)$ , where  $F : n \rightarrow m$  is a morphism in  $\text{Ctx}(\mathcal{R})$  and  $\eta : F \rightarrow F^*(G)$  is a morphism in  $\mathcal{M}^n \rightarrow \mathcal{M}$
- the identity on  $(n, F)$  is the pair  $(\text{id}_n, \text{id}_F)$ , where  $\text{id}_n : n \rightarrow n$  is the identity in  $\text{Ctx}(\mathcal{R})$  and  $\text{id}_F : F \rightarrow F$  is the identity in  $\mathcal{M}^n \rightarrow \mathcal{M}$
- the composition of two morphisms  $(F, \eta_1) : (n, F) \rightarrow (m, G)$  and  $(G, \eta_2) : (m, G) \rightarrow (k, H)$  is the pair  $(G \circ F, F^*(\eta_2) \circ \eta_1)$ , where the first composition is in  $\text{Ctx}(\mathcal{R})$  and the second composition is in  $\mathcal{M}^n \rightarrow \mathcal{M}$

This is a standard (op)Grothendieck construction, and results in a category whose objects can be understood as pairing a kinding context and a typing context over it, and whose morphisms can be understood as simultaneous substitutions.

To appropriately interpret arrow types will we need to know that the category  $\mathcal{M}^n \rightarrow \mathcal{M}$  is cartesian closed. We define:

**Definition 4.6.** A reflexive graph category with isomorphisms  $\mathcal{R}$  has *terminal objects* if each  $X(l)$  has a terminal object  $1_{X(l)}$ . The terminal objects are *stable under face maps* if, for all  $\star \in \mathbf{Bool}$ , the canonical morphism from  $X(\mathbf{f}_\star) 1_{X(1)}$  to  $1_{X(0)}$  is the identity. The terminal objects are *stable under degeneracies* if the canonical morphism from  $X(\mathbf{d}) 1_{X(0)}$  to  $1_{X(1)}$  is in  $\mathcal{M}(1)$ .

**Definition 4.7.** A reflexive graph category with isomorphisms  $\mathcal{R}$  has *products* if each  $X(l)$  has products  $\times_l$  that preserve membership in  $\mathcal{M}(l)$ . The products are *stable under face maps* if, for all  $\star \in \mathbf{Bool}$  and objects  $A, B$  in  $X(1)$ , the canonical morphism from  $X(\mathbf{f}_\star)(A \times_1 B)$  to  $(X(\mathbf{f}_\star)A) \times_0 (X(\mathbf{f}_\star)B)$  is the identity. It has products *stable under degeneracies* if, for all objects  $A, B$  in  $X(0)$ , the canonical morphism from  $X(\mathbf{d})(A \times_0 B)$  to  $(X(\mathbf{d})A) \times_1 (X(\mathbf{d})B)$  is in  $\mathcal{M}(1)$ .

**Definition 4.8.** A reflexive graph category with isomorphisms  $\mathcal{R}$  that has products also has *exponentials* if each  $X(l)$  has exponentials  $\Rightarrow_l$  that preserve membership in  $\mathcal{M}(l)$ . The exponentials are *stable under face maps* if, for all  $\star \in \mathbf{Bool}$  and objects  $A, B$  in  $X(1)$ , the canonical morphism from  $X(\mathbf{f}_\star)(A \Rightarrow_1 B)$  to  $(X(\mathbf{f}_\star)A) \Rightarrow_0 (X(\mathbf{f}_\star)B)$  is the identity. It has exponentials *stable under degeneracies* if, for all objects  $A, B$  in  $X(0)$ , the canonical morphism from  $X(\mathbf{d})(A \Rightarrow_0 B)$  to  $(X(\mathbf{d})A) \Rightarrow_1 (X(\mathbf{d})B)$  is in  $\mathcal{M}(1)$ .

We combine the above to obtain the main definition of this section:

**Definition 4.9.** A reflexive graph category with isomorphisms is *cartesian closed* if it has terminal objects, products, and exponentials, all stable under face maps and degeneracies.

**Example 4.10.** [PER model, continued] Terminal objects, products, and exponentials are defined for  $\mathcal{R}_{PER}$  in the obvious ways, inheriting from the corresponding constructs on PERs. It is not hard to check that all of these constructs are preserved on the nose by the two face maps (projections) and the degeneracy (equality functor), and thus, in our terminology, are stable under face maps and degeneracies.

**Example 4.11.** [Both versions of Reynolds' model, continued] Here, too, terminal objects, products, and exponentials are defined for  $\mathcal{R}_{REY}$  and  $\mathcal{R}_{CREY}$  in the obvious ways, relating two pairs iff their first and second components are related, and two functions iff they map related arguments to related results. It is easy to see that all of these constructs are preserved on the nose (i.e., up to *definitional* equality) by the projections, and thus are stable under face maps. Unlike in the PER model though, they are only preserved by the equality functor  $\text{Eq}$  up to (the canonical) isomorphism. For example, as discussed just after Definition 2.6, the two types  $\text{Id}((a, b), (c, d))$  and  $\text{Id}(a, c) \times \text{Id}(b, d)$  for  $(a, b), (c, d) : A \times B$  are not necessarily identical, although they are isomorphic under the canonical (iso)morphism from  $\text{Eq}(A \times B)$  to  $\text{Eq}(A) \times \text{Eq}(B)$ . A similar situation arises for function types  $A \rightarrow B$ : by function extensionality,  $\text{Id}(f, g)$  and  $\Pi_{a, a':A} \text{Id}(f(a), g(a'))$  are isomorphic, but not necessarily identical, via the canonical isomorphism. Nevertheless, we still get stability under degeneracies since we explicitly allowed for this possibility in Definition 4.8.

## 5 Reflexive Graph Models of Parametricity

As Examples 4.10 and 4.11 show, cartesian closed reflexive graph categories with isomorphisms suitably generalize the structure of sets and relations. Moreover, they allow us to interpret unit, product, and function types in a natural way. To show this, we introduce the following terminology, presented in a form more general than we need for interpreting the simply-typed fragment of System F, but paralleling the later terminology used for interpreting the impredicative fragment.

**Definition 5.1.** A  $\lambda^{\rightarrow}$ -fibration is a split fibration  $U : \mathcal{E} \rightarrow \mathcal{B}$  satisfying the following properties:

1.  $U$  has a split generic object  $\Omega$  in  $\mathcal{B}$ .
2.  $\mathcal{B}$  has a terminal object and products  $(-) \times \Omega$ , and for every object  $I$  in  $\mathcal{B}$ , we have  $I \cong \Omega^n$  for some  $n \in \mathbb{N}$ .
3. Every fiber  $\mathcal{E}_I$  for  $I$  in  $\mathcal{B}$  is cartesian closed, with a terminal object  $1_I$ , products  $\times_I$ , and exponentials  $\Rightarrow_I$ .
4. Beck-Chevalley: for any morphism  $f : I \rightarrow J$  in  $\mathcal{B}$  and objects  $X, Y$  in  $\mathcal{E}_J$ , the canonical morphisms below are isomorphisms:

$$\theta_1(f) : f^*(1_J) \rightarrow 1_I$$

$$\theta_{\times}(f, X, Y) : f^*(X \times_J Y) \rightarrow (f^*(X) \times_I f^*(Y))$$

$$\theta_{\Rightarrow}(f, X, Y) : f^*(X \Rightarrow_J Y) \rightarrow (f^*(X) \Rightarrow_I f^*(Y))$$

A  $\lambda^{\rightarrow}$ -fibration is *split* if these canonical morphisms are identities.

Using a similar idea as in the proof of Lemma 5.2.4 of [7], we can show:

**Lemma 5.2.** *Every  $\lambda^{\rightarrow}$ -fibration is equivalent to a split  $\lambda^{\rightarrow}$ -fibration in a canonical way.*

We now come to our main technical lemma:

**Theorem 5.3.** *Given a cartesian closed reflexive graph category  $\mathcal{R}$  with isomorphisms, the forgetful functor from the category  $\int_n \mathcal{M}^n \rightarrow \mathcal{M}$  to  $\text{Ctx}(\mathcal{R})$  is a split  $\lambda^{\rightarrow}$ -fibration.*

*Proof sketch.* Terminal objects, products, and exponentials in each fiber are given levelwise and pointwise, and hence commute with substitution on the levelwise and pointwise, and hence commute with substitution on the nose. For the full details, we refer the reader to [17].  $\square$

To interpret  $\forall$ -types we need to know that, in the forgetful fibration from Lemma 5.3, each weakening functor induced by the first projection from  $n+1$  to  $n$  for  $n \in \mathbb{N}$  has a right adjoint  $\forall_n$ . Here we differ from [2], where only  $\forall_0$  is required, with the intention that  $\forall_n$  can be derived from  $\forall_0$  using partial application. We observe that this approach does not appear to work since a partial application of an indexed functor is not necessarily an indexed functor. Hence we require an entire family of adjoints  $\forall_n$ .

**Example 5.4** (PER model, continued). Define the adjoint  $\forall_n$  by

$$\begin{aligned} \forall_n \mathcal{F}(0) \bar{A} &:= \{(m, k) \mid \text{for all } A, (m, k) \in \mathcal{F}(0)(\bar{A}, A), \\ &\quad \text{and for all } R, \langle m, k \rangle \in \mathcal{F}(1)(\overline{\text{Eq}} \bar{A}, R)\} \\ \forall_n \mathcal{F}(1) \bar{R} &:= \left( (\forall_n \mathcal{F}(0) \bar{R}_d, \forall_n \mathcal{F}(0) \bar{R}_c), \right. \\ &\quad \left. \{m \mid \text{for all } R, m \in \mathcal{F}(1)(\bar{R}, R)\} \right) \end{aligned}$$

where for any relation  $R := ((A_d, A_c), R_A)$  we write  $R_d$  for  $A_d$  and  $R_c$  for  $A_c$ . We will employ a similar convention for Reynolds' model. To define  $\forall_n$  on a morphism  $\eta : \mathcal{F} \rightarrow \mathcal{G}$ , we put

$$\forall_n \eta(0) \bar{A} := \left( (\forall_n \mathcal{F}(0) \bar{A}, \forall_n \mathcal{G}(0) \bar{A}), \right. \\ \left. \{m \cdot 0\}_{(\forall_n \mathcal{F}(0) \bar{A}) \rightarrow (\forall_n \mathcal{G}(0) \bar{A})} \right)$$

Here,  $m$  is any natural number realizing  $\eta(0) \bar{A}$ . Crucial observations are that all natural transformations are “uniformly realized” in the sense that there is a natural number realizing each such transformation, and since all PERs are defined to be realized by all natural numbers, each is suitably uniform. In particular, if  $\eta$  were not uniformly realized in the above sense then  $\forall_n$  would not be well-defined on morphisms. These observations can be used to show that, in the category-theoretic setting (rather than the setting of  $\omega$ -sets), the family of adjoints  $\forall$  cannot exist precisely because *ad hoc* natural transformations – i.e., natural transformations that are not uniformly realizable, even though each of their components may indeed be realizable – are not excluded.

**Example 5.5** (Reynolds' model, continued). On sets, the adjoint  $\forall_n$  is defined as follows:

$$\forall_n \mathcal{F}(0) \bar{A} := \{f_0 : \Pi_{A:\mathbb{U}} \mathcal{F}(0)(\bar{A}, A) \ \& \\ f_1 : \Pi_{R:R_0} \mathcal{F}(1)(\overline{\text{Eq}} \bar{A}, R)(f_0(R_d), f_0(R_c))\}$$

On relations, we define  $\forall_n \mathcal{F}(1) \bar{R}$  to be the relation with domain  $\forall_n \mathcal{F}(0) \bar{R}_0$  and codomain  $\forall_n \mathcal{F}(0) \bar{R}_1$  mapping  $((f_0, f_1), (g_0, g_1))$  to  $\Pi_{R:R_0} \mathcal{F}(1)(\bar{R}, R)(f_0(R_c), g_0(R_c))$ .

To see that the above definition indeed gives a degeneracy-preserving reflexive graph functor, fix  $\bar{A}$ . We want to show that the two relations  $\text{Eq}(\forall_n \mathcal{F}(0) \bar{A})$  and  $\forall_n \mathcal{F}(1) \overline{\text{Eq}} \bar{A}$  are isomorphic. The domains and codomains of these relations are all the same –  $\forall_n \mathcal{F}(0) \bar{A}$  – so we let both of the underlying maps of the isomorphism be identities (as also required by the coherence condition on the isomorphism and, independently, the definition of a relevant isomorphism). Fix  $((f_0, f_1), (g_0, g_1)) : (\forall_n \mathcal{F}(0) \bar{A}) \times (\forall_n \mathcal{F}(0) \bar{A})$ . We need functions going back and forth between  $\text{Id}((f_0, f_1), (g_0, g_1))$  and  $\Pi_{R:R_0} \mathcal{F}(1)(\overline{\text{Eq}} \bar{A}, R)(f_0(R_0), g_0(R_1))$ . Such functions will automatically be mutually inverse since the types in question are propositions.

Going from left to right is easy using  $\text{Id}$ -induction and  $f_1$ . To go from right to left, fix  $\phi : \Pi_{R:R_0} \mathcal{F}(1)(\overline{\text{Eq}} \bar{A}, R)(f_0(R_0), g_0(R_1))$ . To show  $\text{Id}((f_0, f_1), (g_0, g_1))$  it suffices to show  $\text{Id}(f_0, g_0)$  since the type of  $g_1$  (or  $f_1$ ) is a proposition. By function extensionality, it

suffices to show pointwise equality between  $f_0$  and  $g_0$ . So fix  $B$ . The only thing we can do with  $\phi$  is to apply it to  $\text{Eq}(B)$ , which gives us  $\phi(\text{Eq}(B)) : \mathcal{F}(1)(\overline{\text{Eq}} \bar{A}, \text{Eq}(B))(f_0(B), g_0(B))$ . The relation  $\mathcal{F}(1)(\overline{\text{Eq}} \bar{A}, \text{Eq}(B))$  is isomorphic to  $\text{Eq} \mathcal{F}(0)(\bar{A}, B)$  via  $\varepsilon_{\mathcal{F}(\bar{A}, B)}^{-1}$ . Applying  $\varepsilon_{\mathcal{F}(\bar{A}, B)}^{-1}$  to  $(f_0(B), g_0(B))$  and  $\phi(\text{Eq}(B))$  thus gives us  $\text{Id}(\varepsilon_{\mathcal{F}(\bar{A}, B)}^{-1} f_0(A), \varepsilon_{\mathcal{F}(\bar{A}, B)}^{-1} g_0(B))$ . The coherence condition on  $\varepsilon_{\mathcal{F}}$  tells us that the respective images  $\varepsilon_{\mathcal{F}(\bar{A}, B)\top}$  and  $\varepsilon_{\mathcal{F}(\bar{A}, B)\perp}$  of  $\varepsilon_{\mathcal{F}(\bar{A}, B)}$  under the two face maps are the identity on  $\mathcal{F}(0)(\bar{A}, B)$ , and thus are  $\varepsilon_{\mathcal{F}(\bar{A}, B)\top}^{-1}$  and  $\varepsilon_{\mathcal{F}(\bar{A}, B)\perp}^{-1}$ . This gives  $\text{Id}(f_0(A), g_0(B))$  as desired.

**Example 5.6** (A categorical version of Reynolds' model, continued). On sets, the adjoint  $\forall_n$  is defined as follows:

$$\begin{aligned} \forall_n \mathcal{F}(0) \bar{A} &:= \{f_0 : \Pi_{A:\mathbb{U}} \mathcal{F}(0)(\bar{A}, A) \ \& \\ &\quad f_1 : \Pi_{R:R_0} \mathcal{F}(1)(\overline{\text{Eq}} \bar{A}, R)(f_0(R_d), f_0(R_c)) \ \& \\ &\quad \Pi_{i:M(0)} \mathcal{F}(0)(\text{id}_{M(0)}(A), i) f_0(i_d) = f_0(i_c)\} \end{aligned}$$

The last condition says that  $f_0$  is functorial in its argument, in the sense that if  $i$  is an isomorphism between two types  $A, B : \text{Set}_0$ , then  $f_0(A)$  and  $f_0(B)$  are suitably related via the obvious isomorphism between  $\mathcal{F}(0)(\bar{A}, A)$  and  $\mathcal{F}(0)(\bar{A}, B)$ . This condition, which does not have an analogue in the set-theoretic presentation of Reynolds' model, is needed because we do not work with discrete domains (e.g., we use  $\mathcal{F} : \mathcal{M}^n \rightarrow \mathcal{M}$  rather than  $\mathcal{F} : |\mathcal{M}|^n \rightarrow \mathcal{M}$ ), as is common in other presentations of parametricity. A very similar condition does appear, e.g., in the definition of parametric limits for the category of sets in [2]. The analogous condition asserting the functoriality of  $f_1$  is automatically satisfied since the codomain of  $f_1$  is a proposition. On relations, we use the same definition as in Example 5.5.

**Definition 5.7.** A  $\lambda 2$ -fibration is a  $\lambda^\top$ -fibration  $U : \mathcal{E} \rightarrow \mathcal{B}$  satisfying the following properties:

1. For each  $I$  in  $\mathcal{B}$ , the weakening functor induced by the first projection from  $I \times \Omega$  to  $\Omega$  has a right adjoint  $\forall_I$ .
2. Beck-Chevalley: for any morphism  $f : I \rightarrow J$  in  $\mathcal{B}$  and object  $X$  in  $\mathcal{E}_J$ , the canonical morphism below is an isomorphism:

$$\theta_{\forall}(f, X) : f^*(\forall_J(X)) \rightarrow \forall_I((f \times \text{id})^*(X))$$

A  $\lambda 2$ -fibration is *split* if it is a split  $\lambda^\top$ -fibration and the canonical morphism above is the identity.

Seely [16] essentially showed the following:

**Theorem 5.8** (Seely). *Every split  $\lambda 2$ -fibration  $U : \mathcal{E} \rightarrow \mathcal{B}$  gives a sound model of System F in which:*

- every type context  $\Gamma$  is interpreted as an object  $\llbracket \Gamma \rrbracket$  in  $\mathcal{B}$
- every type  $\Gamma \vdash T$  is interpreted as an object  $\llbracket \Gamma \vdash T \rrbracket$  in the fiber over  $\llbracket \Gamma \rrbracket$
- every term context  $\Gamma; \Delta$  is interpreted as an object  $\llbracket \Gamma \vdash \Delta \rrbracket$  in the fiber over  $\llbracket \Gamma \rrbracket$
- every term  $\Gamma; \Delta \vdash t : T$  is interpreted as a morphism  $\llbracket \Gamma; \Delta \vdash t : T \rrbracket$  from  $\llbracket \Gamma; \Delta \rrbracket$  to  $\llbracket \Gamma \vdash T \rrbracket$  in the fiber over  $\llbracket \Gamma \rrbracket$

A (not necessarily split)  $\lambda 2$ -fibration also gives a sound model of System F, due to the following:

**Lemma 5.9.** *Every  $\lambda 2$ -fibration is equivalent to a split  $\lambda 2$ -fibration in a canonical way.*

We now want to specify when a model of System F given by a  $\lambda 2$ -fibration is relationally parametric. If  $\mathcal{R}$  is a cartesian closed



reflexive graph category with isomorphisms, we denote by  $F(\mathcal{R})$  the  $\lambda^{\rightarrow}$ -fibration induced by  $\mathcal{R}$  as in Theorem 5.3. To formulate an abstract definition of a parametric model, we will appropriately relate a  $\lambda 2$ -fibration  $U$  to  $F(\mathcal{R})$ . To see how, we revisit the simplest model, namely the System F term model. In the  $\lambda 2$ -fibration  $U_{term}$  corresponding to the term model, the fiber over  $n \in \mathbb{N}$  consists of types and terms with  $n$  free type variables. Let  $\mathcal{U}$  be the category consisting of closed System F types and terms between them. Then  $\mathcal{U}$  induces a  $\lambda^{\rightarrow}$ -fibration,  $U_{set}$ , whose fiber over  $n$  consists of functors  $|\mathcal{U}|^n \rightarrow \mathcal{U}$  and natural transformations between them.

A type  $\bar{\alpha} \vdash T$  with  $n$  free variables can now be seen as functor  $|\mathcal{U}|^n \rightarrow \mathcal{U}$ , and a term  $\bar{\alpha}; x : S \vdash t : T$  as a natural transformation between  $S$  and  $T$ . We thus have a morphism of  $\lambda^{\rightarrow}$ -fibrations  $\mu : U_{term} \rightarrow U_{set}$ . However, unlike  $U_{term}$ ,  $U_{set}$  does not admit the family of adjoints required to make it a  $\lambda 2$ -fibration. Still, we can view  $U_{term}$  as a version of  $U_{set}$  that “enriches” the functors and natural transformations with enough extra information to ensure that the desired adjoints exist: in this example, the information that the maps involved are not *ad hoc*, but come from syntax. Since these adjunctions are only applicable to non-empty contexts, no such “enrichment” should be necessary for objects and morphisms over the terminal object. And indeed, the restriction of  $\mu$  to the fibers over the respective terminal objects is clearly an equivalence. These observations echo those immediately following Definition 2.1, and motivate our main definition:

**Definition 5.10.** Let  $\mathcal{R}$  be a cartesian closed reflexive graph category with isomorphisms. A *parametric model of System F over  $\mathcal{R}$*  is a  $\lambda 2$ -fibration  $U$  together with a morphism  $\mu : U \rightarrow F(\mathcal{R})$  of  $\lambda^{\rightarrow}$ -fibrations whose restriction to the fibers of  $U$  and  $F(\mathcal{R})$  over the terminal objects is full, faithful, and essentially surjective.

Our main theorem shows that the definition of a parametric model is indeed sensible:

**Theorem 5.11.** *Every parametric model of System F over a cartesian closed reflexive graph category  $(\mathcal{X}, (\mathcal{M}, \mathcal{I}))$  with isomorphisms, as specified in Definition 5.10, is a sound model in which:*

- every type  $\Gamma \vdash T$  can be seen as a face map- and degeneracy-preserving reflexive graph functor  $\llbracket \Gamma \vdash T \rrbracket : \mathcal{M}^{|\Gamma|} \rightarrow \mathcal{M}$
- every term  $\Gamma; \Delta \vdash t : T$  can be seen as a face map- and degeneracy-preserving reflexive graph natural transformation  $\llbracket \Gamma; \Delta \vdash t : T \rrbracket : \llbracket \Gamma \vdash \Delta \rrbracket \rightarrow \llbracket \Gamma \vdash T \rrbracket$ , with the domain and codomain seen as reflexive graph functors into  $\mathcal{X}$

**Theorem 5.12** (PER model). *Let  $\mathcal{R}_{PER}$  be the cartesian closed reflexive graph category with isomorphisms defined in Examples 2.3, 3.3, and 4.10. The family of adjoints defined in Example 5.4 makes  $F(\mathcal{R}_{PER})$  into a  $\lambda 2$ -fibration, and hence into a parametric model of System F over  $\mathcal{R}_{PER}$ .*

**Theorem 5.13** (Reynolds’ model). *Let  $\mathcal{R}_{REY}$  be the reflexive graph category with isomorphisms defined in Examples 2.4, 3.5, and 4.11. The family of adjoints defined in Example 5.5 makes  $F(\mathcal{R}_{REY})$  into a  $\lambda 2$ -fibration, and hence into a parametric model of System F over  $\mathcal{R}_{REY}$ .*

**Theorem 5.14** (A categorical version of Reynolds’ model). *Let  $\mathcal{R}_{CREY}$  be the reflexive graph category with isomorphisms defined in Examples 2.4, 3.4, and 4.11. The family of adjoints defined in Example 5.6 makes  $F(\mathcal{R}_{CREY})$  into a  $\lambda 2$ -fibration, and hence into a parametric model of System F over  $\mathcal{R}_{CREY}$ .*

## 6 A Proof-Relevant Model of Parametricity

We now describe a proof-relevant version of Reynolds’ model, in which witnesses of relatedness are themselves related. The construction of such a model is the subject of [12], but the development there seems to contain a major technical gap. Specifically, it is unclear how to prove the  $\forall$ -case in Lemma 9.4: when types are interpreted as discrete functors  $|\mathcal{X}|^n \rightarrow \mathcal{X}$ , the reindexing of a degeneracy-preserving functor might not be degeneracy-preserving. We already observed this in the introduction, but this issue is not addressed in [12] and the proof of the lemma is not given there. Since this lemma is crucial to the soundness of the interpretation, it is unknown whether the result of [12] can be salvaged as-is. For this reason, we only reuse the main ideas of [12] for handling the higher dimensional structure and otherwise proceed independently.

**Example 6.1.** We use the same ambient category as in Example 2.4 and reuse the (internal) category Set of types. The category  $\mathbf{R}$  of relations is almost the same as in Example 2.4, except that relations are now proof-relevant, i.e.,  $R_0 := \Sigma_{A,B:\text{Set}} A \times B \rightarrow \mathbb{U}$ . Given relations  $R$  on  $A, B$  and  $S$  on  $C, D$ , to relate two witnesses  $p : R(a, b)$  and  $q : S(c, d)$  we should know *a priori* how  $a$  relates to  $c$  and  $b$  to  $d$ . This motivates defining the category  $2\mathbf{R}$  of *2-relations*, whose objects  $Q$  are tuples  $(Q^0, Q^1, Q^2, Q^3)$  of relations forming a square

$$\begin{array}{ccc} A & \xrightarrow{Q^0} & B \\ Q^1 \downarrow & & \downarrow Q^3 \\ C & \xrightarrow{Q^2} & D \end{array}$$

together with a Prop-valued predicate (also denoted  $Q$ ) on the type of tuples of the form  $((a, b, c, d), (p, q, r, s))$ , where  $p : Q^0(a, b)$ ,  $q : Q^1(a, c)$ ,  $r : Q^2(c, d)$ , and  $s : Q^3(b, d)$ . This gives four face maps from  $2\mathbf{R}$  to  $\mathbf{R}$ , one for each edge. We also have four functors in the other direction: e.g., given  $R$ , we obtain the 2-relation  $\text{Eq}_=(R)$  by placing  $R$  on top and bottom, with equalities as vertical edges, and mapping  $((a, b, a, b), (p, -, r, -))$  to  $\text{Id}(p, r)$ . Similarly,  $\text{Eq}_\parallel(R)$  places  $R$  on left and right,  $\text{C}_\top(R)$  places  $R$  on top and left, and  $\text{C}_\perp(R)$  places  $R$  on bottom and right, all filling the remaining edges with equalities. The functors  $\text{Eq}_=$ ,  $\text{Eq}_\parallel$  are called *degeneracies* and  $\text{C}_\top$ ,  $\text{C}_\perp$  are called *connections*. We define terminal objects, products, exponentials, and isomorphisms in the obvious way.

The above structure induces two  $\lambda^{\rightarrow}$ -fibrations of interest: the first one is induced by combining the first two levels, the categories Set and  $\mathbf{R}$ , into a cartesian closed reflexive graph category with isomorphisms  $\mathbf{R}_{PREY}$ ; this is the fibration  $F(\mathbf{R}_{PREY})$ . We recall that the objects of  $F(\mathbf{R}_{PREY})$  over  $n$  are pairs  $\{\mathcal{F}(I) : \mathcal{M}(I)^n \rightarrow \mathcal{M}(I)\}_{I \in \{0,1\}}$  of functors that commute with the two face maps from  $\mathbf{R}$  to Set on the nose, as well as with the degeneracy  $\text{Eq}$  up to a suitably coherent natural isomorphism. The morphisms are pairs  $\{\eta(I) : \mathcal{F}(I) \rightarrow \mathcal{G}(I)\}_{I \in \{0,1\}}$  of natural transformations that respect both face maps from  $\mathbf{R}$  to Set and the degeneracy  $\text{Eq}$ . The second fibration, which we call  $\mathbf{F}_{2D}$ , is induced in much the same way, but taking into account all three levels. This means that the objects over  $n$  are triples  $\{\mathcal{F}(I) : \mathcal{M}(I)^n \rightarrow \mathcal{M}(I)\}_{I \in \{0,1,2\}}$  of functors that commute with *all* face maps — the two from from  $\mathbf{R}$  to Set as well as the four from  $2\mathbf{R}$  to  $\mathbf{R}$  — on the nose and *all* degeneracies  $\text{Eq}$ ,  $\text{Eq}_=$ ,  $\text{Eq}_\parallel$  and connections  $\text{C}_\top$ ,  $\text{C}_\perp$  up to suitably coherent natural isomorphisms. Analogously, the morphisms are triples  $\{\eta(I) : \mathcal{F}(I) \rightarrow \mathcal{G}(I)\}_{I \in \{0,1,2\}}$  of natural transformations that respect all face maps, degeneracies, and connections. We have the obvious forgetful morphism of  $\lambda^{\rightarrow}$ -fibrations

from  $F_{2D}$  to  $F(R_{PREY})$  that only retains the structure pertaining to levels 0 and 1.

The fibration  $F_{2D}$  admits a family of adjoints to weakening functors as follows. The adjoint  $\forall_n \mathcal{F}(0) \bar{A}$  is the type

$$\begin{aligned} & \{f_0 : \Pi_{A:\mathbb{U}} \mathcal{F}(0)(\bar{A}, A) \& \\ & f_1 : \Pi_{R:R_0} \mathcal{F}(1)(\overline{\text{Eq}} \bar{A}, R) (f_0(R_d), f_0(R_d)) \& \\ & f_2 : \Pi_{Q:2R_0} \mathcal{F}(2)(\overline{\text{Eq}}_=(\text{Eq}(A)), Q) ((f_0 Q_d^0, f_0 Q_c^0, f_0 Q_c^1, f_0 Q_c^2), \\ & \quad (f_1 Q^0, f_1 Q^1, f_1 Q^2, f_1 Q^3)) \& \\ & \Pi_{i:M(0)} \mathcal{F}(0)(\overline{\text{id}}_{M(0)}(A), i) f_0(i_c) = f_0(i_d) \& \\ & \Pi_{i:M(1)} \mathcal{F}(1)(\overline{\text{id}}_{M(1)}(\text{Eq} A), i) (f_0(i_d)_d, f_0(i_d)_c) f_1(i_d) = f_1(i_c) \} \end{aligned}$$

In the type of  $f_2$ , we could have just as well used any of the other functors  $\text{Eq}_{\parallel}$ ,  $\text{C}_{\top}$ ,  $\text{C}_{\perp}$  instead of  $\text{Eq}_=$  since their compositions with  $\text{Eq}$  are all naturally isomorphic. We next define  $\forall_n \mathcal{F}(1) \bar{R}$  to be the relation with domain  $\forall_n \mathcal{F}(0) \bar{R}_d$  and codomain  $\forall_n \mathcal{F}(0) \bar{R}_c$  mapping  $((f_0, f_1, f_2), (g_0, g_1, g_2))$  to

$$\begin{aligned} & \{\phi : \Pi_{R:R_0} \mathcal{F}(1)(\bar{R}, R) (f_0(R_d), g_0(R_c)) \& \\ & \phi_=: \Pi_{Q:2R_0} \mathcal{F}(2)(\overline{\text{Eq}}_=\bar{R}, Q) ((f_0 Q_d^0, f_0 Q_c^0, g_0 Q_c^1, g_0 Q_c^2), \\ & \quad (f_1 Q^0, \phi Q^1, g_1 Q^2, \phi Q^3)) \& \\ & \dots \& \\ & \Pi_{i:M(1)} \mathcal{F}(1)(\overline{\text{id}}_{M(1)}(R), i) (f_0(i_d)_d, g_0(i_d)_c) \phi_1(i_d) = \phi_1(i_c) \} \end{aligned}$$

The component  $\phi_=:$  asserts that  $\phi$  appropriately interacts with the degeneracy  $\text{Eq}_=$ . The analogous components  $\phi_{\parallel}$ ,  $\phi_{\text{C}_{\top}}$ ,  $\phi_{\text{C}_{\perp}}$  for  $\text{Eq}_{\parallel}$ ,  $\text{C}_{\top}$ ,  $\text{C}_{\perp}$  are omitted for space reasons.

We define  $\forall_n \mathcal{F}(2) \bar{Q}$  to be the 2-relation with underlying tuple of relations  $(\forall_n \mathcal{F}(1) \bar{Q}^0, \forall_n \mathcal{F}(1) \bar{Q}^1, \forall_n \mathcal{F}(1) \bar{Q}^2, \forall_n \mathcal{F}(1) \bar{Q}^3)$ , and mapping a tuple of the form  $((f_0, \dots), (g_0, \dots), (h_0, \dots), (l_0, \dots), ((\phi_0, \dots), (\phi_1, \dots), (\phi_2, \dots), (\phi_3, \dots)))$  to

$$\Pi_{Q:2R_0} \mathcal{F}(2)(\bar{Q}, Q) ((f_0 Q_d^0, g_0 Q_c^0, h_0 Q_d^1, l_0 Q_c^2), (f_0 Q^0, \phi_1 Q^1, \phi_2 Q^2, \phi_3 Q^3))$$

Finally, unlike the frameworks [2, 4, 5, 7, 10, 15], our definition of a parametric model recognizes the above proof-relevant model:

**Theorem 6.2** (Proof-relevant model). *The family of adjoints defined in Example 6.1 makes  $F_{2D}$  into a  $\lambda_2$ -fibration, and hence into a parametric model of System F over  $R_{PREY}$ .*

## 7 Discussion

We can now be more specific about how our approach compares to the external approaches in [4, 7, 10, 15], all of which are based on a reflexive graph of  $\lambda_2$ -fibrations. The definition in [4] appears to be too restrictive: it requires a comprehension structure that, e.g., the  $\lambda_2$ -fibration corresponding to Reynolds' model does not admit. In addition, none of these frameworks seem to recognize the  $\lambda$ -fibration corresponding to the proof-relevant model as parametric; indeed, it is unclear how to define the family of adjoints for the second fibration (called  $r$  in [7]) of "heterogeneous" reflexive graph functors in a way that is compatible with the adjoint structure on the original  $\lambda_2$ -fibration. This is because, unlike in the proof-irrelevant case, the definition of  $\forall_n \mathcal{F}(1)$  now has conditions, such as the one witnessed by  $\phi_=:$ , which are only meaningful for "homogeneous" reflexive graph functors, i.e., those where the domain and codomain of  $\mathcal{F}(1)(\bar{R})$  are given by the same functor  $\mathcal{F}(1)$ , albeit applied to

different arguments ( $\overline{R}_d$  vs.  $\overline{R}_c$ ). Our definition does not rely on or require two compatible adjoint structures, which is why we are indeed able to recognize the proof-relevant model as parametric.

We indicate three directions for future work. Readers interested in *applications of parametricity* will notice that we do not require conditions such as fullness or (op)cartesianness of certain maps or well-pointedness of certain categories. This follows the spirit of [7], where the notion of *parametricity* pertains to the suitable interaction with (what we call) face maps and degeneracies. Specific *applications* such as establishing the Graph Lemma and the existence of initial algebras are left for another occasion. Readers fond of *type theory* might wonder about possible models expressed in the *intensional* version of dependent type theory. Although currently there are no well-known models for which the latter would be the right choice of meta-theory, that might change with more research into higher notions of parametricity. Finally, readers familiar with *cubical sets* no doubt recognized the structure of sets, relations, and 2-relations with face maps, degeneracies, and connections from the last section as the first few levels of the cubical hierarchy, and wonder whether one can formulate the analogous notion of 2-parametricity, 3-parametricity, . . . using this hierarchy. We conjecture the answer to be a *YES!* and plan to pursue this question in future work.

**Acknowledgments** This research is supported by NSF awards 1420175 and 1545197. We thank Steve Awodey and Peter Dybjer for helpful discussions. We also thank the anonymous referee who independently suggested the formulation of Reynolds' model in Example 6.1, which appeared in our earlier preprint [8].

## References

- [1] R. Atkey, N. Ghani, and Patricia Johann. 2014. A Relationally Parametric Model of Dependent Type Theory. In *Principles of Programming Languages*. 503–516.
- [2] B. Dunphy and U. Reddy. 2004. Parametric Limits. In *Logic in Computer Science*. 242–251.
- [3] N. Ghani, F. Nordvall Forsberg, and F. Orsanigo. 2016. Proof-Relevant Parametricity. In *A List of Successes That Can Change the World - Essays Dedicated to Philip Wadler on the Occasion of His 60th Birthday (Lecture Notes in Computer Science 9600)*. 109–131.
- [4] N. Ghani, F. Nordvall Forsberg, and A. Simpson. 2016. Comprehensive Parametric Polymorphism: Categorical Models and Type Theory. In *Foundations of Software Science and Computation Structures*. 3–19.
- [5] N. Ghani, P. Johann, F. Nordvall Forsberg, F. Orsanigo, and T. Revell. 2015. Bifibrational Functorial Semantics for Parametric Polymorphism. In *Mathematical Foundations of Program Semantics*. 165–181.
- [6] J.-Y. Girard, P. Taylor, and Y. Lafont. 1989. *Proofs and Types*. Cambridge University Press.
- [7] B. Jacobs. 1999. *Categorical Logic and Type Theory*. Elsevier.
- [8] P. Johann and K. Sojakova. 2017. Cubical Categories for Higher-Dimensional Parametricity. (2017). Available on arxiv.org as arXiv:1701.06244.
- [9] G. Longo and E. Moggi. 2009. Constructive Natural Deduction and its  $\omega$ -set Interpretation. *Mathematical Structures in Computer Science* 1(2) (2009), 215–254.
- [10] Q. Ma and J. Reynolds. 1992. Types, Abstraction, and Parametric Polymorphism. In *Mathematical Foundations of Program Semantics*. 1–40.
- [11] P. W. O'Hearn and R. D. Tennent. 1995. Parametricity and Local Variables. *J. ACM* 42(3) (1995), 658–709.
- [12] F. Orsanigo. 2017. *Bifibrational Parametricity: From Zero to Two Dimensions*. Ph.D. Dissertation. University of Strathclyde.
- [13] J. Reynolds. 1983. Types, Abstraction, and Parametric Polymorphism. In *International Federation for Information Processing*. 513–523.
- [14] J. C. Reynolds. 1984. Polymorphism is not set-theoretic. *Semantics of Data Types* (1984), 145–156.
- [15] E. Robinson and G. Rosolini. 1994. Reflexive Graphs and Parametric Polymorphism. In *Logic in Computer Science*. 364–371.
- [16] R. A. G. Seely. 1987. Categorical Semantics for Higher Order Polymorphic Lambda Calculus. *Journal of Symbolic Logic* 52 (1987), 969–989.
- [17] K. Sojakova and P. Johann. 2018. A General Framework for Relational Parametricity. (2018). Available from the last author's webpage.
- [18] P. Wadler. 1987. Theorems for Free!. In *Functional Programming and Computer Architecture*. 347–359.