

# A State Is Its Relations: The Yoneda Lemma and Relational Identity in International Relations

Robert J. Carroll\*

March 9, 2026

## Abstract

This paper proves that, within any theory that specifies actors and relations, a state's identity is nothing more and nothing less than its complete relational profile—the totality of relationships that every other actor bears to it. The proof uses the Yoneda lemma from category theory. I formalize IR theories as categories, models as functors, and translations between models as natural transformations. The Yoneda lemma establishes that every datum a model assigns to a state is a global structural commitment—a coherent reading of the state's position in the relational network, compressed into a single element—and that two states are structurally identical if and only if their relational profiles agree. The framework yields concrete results: applied to two well-known bargaining models of war, the naturality condition identifies exactly which parameter configurations in one model can be microfounded by the other; applied to existing datasets, it reveals that attribute tables encode a theory in which the Yoneda lemma is vacuous. The constructivist claim that identities are constituted by interaction is the correct structural consequence of taking relations seriously: not a philosophical stance but a theorem.

**Keywords:** category theory, Yoneda lemma, relational ontology, state identity, constructivism, metatheory

---

\*Department of Political Science, University of Illinois at Urbana-Champaign. Email: [rjc@illinois.edu](mailto:rjc@illinois.edu). Working paper; comments welcome.

# 1 Introduction

What is the United States? A rationalist might answer: a state with fixed preferences, a given military capability, a GDP of roughly twenty-eight trillion dollars. A constructivist might answer: an actor whose identity is constituted through its interactions with other actors—through trade, alliance, rivalry, recognition. The two answers reflect a deep disagreement about the ontology of international relations, one that has structured the discipline’s theoretical debates for decades (Wendt, 1999; Waltz, 1979; Jackson, 2011).

This paper argues that the constructivist answer is not merely a philosophical preference. It is a theorem.

The argument draws on category theory, a branch of mathematics developed in the 1940s to study the structure of relationships between mathematical objects (Mac Lane, 1998). Category theory is unusual among mathematical frameworks in that it treats relationships—not objects, not properties, not sets—as the primary data. Objects are characterized entirely by how they relate to other objects. This makes it a natural, if unexpected, formal language for a discipline that has spent decades debating whether states are defined by their attributes or by their relations.

The central result is the *Yoneda lemma*, a theorem that holds in any category. It says: an object is completely determined, up to isomorphism, by the totality of relationships that all other objects bear to it. There are no hidden properties, no intrinsic essence, no residual identity lurking behind the relations. Applied to international relations, this yields a precise claim: a state’s identity, within any relational theory, is nothing more and nothing less than its complete relational profile.

## 1.1 What the paper does

The paper develops this argument in several steps, building up the necessary category-theoretic vocabulary from scratch within the context of IR.

A *theory* of international relations is formalized as a category: a collection of actors (objects) connected by typed relations (morphisms) that compose associatively (Section 2). The trade category has states as objects and trade relationships as morphisms; the alliance category has alliance commitments; the recognition category has acts of diplomatic recognition. Each is a different theory of the same system of actors, foregrounding a different type of relation.

A *model* of a theory is a functor—a structure-preserving map from the category to the world of sets and functions (Section 3). The liberal model of the trade category assigns welfare outcomes to states and welfare-transmission functions to trade links. The realist model assigns power positions and power-transmission functions. Both are models of the same theory; they agree on the relational structure and disagree on what it means. Existing datasets—the Penn World Tables, the COW bilateral trade data, the CINC score—are shown to be implicit functors, and the shape of each dataset encodes theoretical commitments that are usually left unexamined.

A *translation* between models is a natural transformation: a systematic, structure-respecting conversion of one model’s data into another’s (Section 4). The naturality condition—the requirement that translation commutes with the relational structure—is demanding, and its failure is informative: two models that admit no natural transformation between them are incommensurable in a precise structural sense. Applied to two well-known bargaining models of war, the naturality condition identifies exactly which parameter configurations in one model can be microfounded by the other, and which cannot.

The *relational profile* of a state is the representable functor  $\text{Hom}(-, X)$ : the complete record of how every other actor relates to  $X$ , with no interpretive overlay (Section 5). This is the maximally innocent model—the theory talking about itself. The *Yoneda lemma* then proves that every model’s assessment of  $X$  is a projection of this relational profile, and the *Yoneda embedding* proves that two states are structurally identical if and only if their profiles agree.

Finally, functors between categories formalize *theory comparison*: the question “does trade cause peace?” becomes a question about the existence of a specific functor from the trade category to a conflict category, subject to compositional constraints (Section 6).

## 1.2 What the paper contributes

The paper makes three contributions.

First, it provides a *formal framework for relational ontology* in IR. The claim that identity is relational has been a staple of constructivist theory since at least [Wendt \(1992\)](#), but it has remained a philosophical stance rather than a formal result. The Yoneda lemma makes it a theorem: any framework that takes relational structure seriously and satisfies the minimal axioms of

a category will, as a matter of mathematical necessity, yield the conclusion that identity is relational profile.

Second, it provides a *precise metatheory* for comparing IR models. The question of whether two models of the same theory are commensurable—whether there is a coherent, structure-respecting translation between them—has a determinate answer in terms of natural transformations. The question of whether two theories with different primitives can be related has a determinate answer in terms of functors between categories. These are not verbal assessments of similarity; they are structural properties with definite values.

Third, it provides *concrete analytical tools*. The naturality condition on two well-known bargaining models of war—Fearon’s exogenous-parameter framework and the Beviá–Corchón contest model—yields a provable constraint that identifies exactly which parameter configurations in one model can be microfounded by the other. The analysis of datasets as functors reveals what is lost when relational structure is flattened into attribute tables. The formalization of “does trade cause peace?” as a functor between categories clarifies what kind of object the question is about.

### 1.3 What the paper does not do

The framework is structural, not dynamic: it captures the relational architecture of the international system at a given moment but does not model change over time. It says nothing about agency, intentionality, or the mechanisms by which relationships are formed and dissolved. And it makes identity theory-dependent: a state’s identity in the trade category may differ from its identity in the alliance category, because the two categories carry different relational structures. These limitations are discussed in Section 7.2, along with extensions—enriched categories, higher categories, dependent type theory—that may address some of them.

### 1.4 A note on mathematics

The paper is self-contained. Every category-theoretic concept used in the argument is defined from scratch, motivated by IR examples, and illustrated with a running four-state trade network (Laos, Vietnam, China, the United States) that threads through every section. The reader needs no prior acquaintance with category theory. The mathematics is elementary—the Yoneda lemma’s proof is a few lines—but its consequences are not. The

goal is to make those consequences available to the IR community, in a form that does not require a second reading to penetrate.

## 1.5 Roadmap

Section 2 introduces categories as theories. Section 3 introduces functors as models. Section 4 introduces natural transformations as translations between models. Section 5 states and interprets the Yoneda lemma. Section 6 develops functors between categories as tools for theory comparison. Section 7 discusses the framework’s relationship to existing IR metatheory, its limitations, and possible extensions. Section 8 concludes.

## 2 What Is a Theory of International Relations?

What does it mean to have a *theory* of international relations? At minimum, a theory specifies a domain of actors and a type of relationship between them. Realism foregrounds power relations among states. Liberalism foregrounds institutional and economic ties. Constructivism foregrounds the social relationships—recognition, enmity, friendship—through which actors constitute one another. In each case, the theory tells us who the relevant actors are and what kind of connection between them matters.

Category theory offers a way to make this intuition precise. A *category* is a mathematical structure that captures exactly this: a collection of objects and a specified type of relationship (called *morphisms*) between them, subject to minimal axioms that any reasonable notion of “relationship” should satisfy. The claim of this paper is that a theory of international relations, understood at the appropriate level of abstraction, *is* a category.

### 2.1 Categories as theories

**Definition 2.1** (Category). A *category*  $\mathcal{C}$  consists of:

- (i) a collection of *objects*,  $\text{Ob}(\mathcal{C})$ ;
- (ii) for each pair of objects  $A, B$ , a set of *morphisms*  $\text{Hom}_{\mathcal{C}}(A, B)$ ;
- (iii) for each triple of objects  $A, B, C$ , a *composition* operation

$$\circ: \text{Hom}(B, C) \times \text{Hom}(A, B) \rightarrow \text{Hom}(A, C);$$

(iv) for each object  $A$ , an *identity morphism*  $\text{id}_A \in \text{Hom}(A, A)$ ;

subject to the axioms:

- **Associativity.** For morphisms  $f: A \rightarrow B$ ,  $g: B \rightarrow C$ ,  $h: C \rightarrow D$ , we have  $h \circ (g \circ f) = (h \circ g) \circ f$ .
- **Identity.** For any morphism  $f: A \rightarrow B$ , we have  $f \circ \text{id}_A = f = \text{id}_B \circ f$ .

The definition is spare, and deliberately so. Each component carries a substantive interpretation in international relations:

**Objects as actors.** The objects of the category are the actors that the theory takes as primary. These are typically states, but nothing in the formalism prevents us from taking non-state actors, international organizations, firms, or individuals as objects when the theory demands it. What matters is that the theory specifies a determinate domain.

**Morphisms as relations.** For each pair of actors  $A$  and  $B$ , the set hom-set  $\text{Hom}(A, B)$  collects the relationships of the relevant type that run from  $A$  to  $B$ . Crucially, morphisms are *typed* and *directed*. A trade relationship is not an alliance relationship; an export flow from Germany to Poland is not the same as an export flow from Poland to Germany. The choice of what counts as a morphism is the fundamental theoretical commitment: it is what distinguishes one theory from another.

But a single category captures only one type of relationship at a time. International politics, of course, involves many types simultaneously—trade, alliance, recognition, rivalry—and much of the interesting structure lives in how they interact. Category theory handles this naturally. Multiple categories on the same set of actors can coexist, and *functors* between them (Section 6) formalize how one relational structure maps onto another. Richer extensions are also available: *higher categories* allow relationships between relationships (a renegotiation of a trade agreement is a morphism between morphisms), and *multi-sorted* or *enriched* categories can encode relationships among more than two actors at once, or weight relationships with quantitative data rather than treating them as present-or-absent. We will not need these extensions in what follows, but it is worth noting that the framework scales: the basic category-theoretic vocabulary introduced here is the entry point to a much larger toolkit, not the whole of it.

**Composition as transitivity.** If  $A$  bears a relationship to  $B$  and  $B$  bears a relationship to  $C$ , then there exists a composite relationship from  $A$  to  $C$ . This is not a trivial assumption. It says that the relational type under study is *transitive* in a structured way: indirect connections, mediated through third parties, are themselves relationships of the same type. When the United States exports to Vietnam and Vietnam exports to Laos, the composition axiom insists that there is a determinate (if indirect) trade relationship from the United States to Laos. The theory must be able to say what it is.

**Identity as self-relation.** Every actor bears a distinguished relationship to itself. In a trade category, this is the domestic economy—the baseline of self-exchange against which foreign trade is measured. In an alliance category, it is the commitment to self-defense that exists prior to any alliance. These are not trivial glosses: any theory of international politics must announce what self-relation consists in, and different theories give different answers. The identity morphism is the formal expression of this commitment.<sup>1</sup>

Notice, though, that some theories endow the identity morphism with rich content (the domestic economy, the commitment to self-defense), while a *purely relational* theory might leave it bare—the uninterpreted sentence “ $A$  is  $A$ ,” carrying no intrinsic content at all. In such a theory, the identity morphism says only that the actor exists and can enter into relations; everything substantive about the actor comes from its connections to others. Once we choose a model  $F$  (Section 3), the identity morphism maps to the identity function on the set  $F(A)$ —a tautology whose content depends entirely on what  $F$  assigns to  $A$ . Whether that assignment is thick or thin is a theoretical choice, and one of the things the Yoneda lemma will clarify.

**Associativity as coherence.** When chaining three or more relationships together, the order of composition does not matter. A chain of trade linkages from  $A$  through  $B$  through  $C$  to  $D$  yields the same composite relationship

---

<sup>1</sup>The word “identity” is doing double duty in this paper. The *identity morphism*  $\text{id}_A$  is a structural element of the category—it says that  $A$  stands in a self-relation and that this self-relation is neutral under composition. The *identity of a state*—the question of what makes  $A$  the actor it is—is the subject of Section 5, where it will turn out to be determined by the totality of  $A$ ’s relationships to every other actor, not by the self-loop alone. The two notions are related but distinct: the identity morphism tells us what a theory takes self-relation to be; the Yoneda lemma tells us that a state’s identity is constituted by all its relations, including but not limited to the self-relation.

regardless of whether we first compose the  $A$ – $B$  and  $B$ – $C$  links or the  $B$ – $C$  and  $C$ – $D$  links. This is a coherence condition: it ensures that the theory’s relational structure is internally consistent, that there is no ambiguity in what it means to trace a path through the network of relationships.

*Remark 2.2.* The reader may notice that a category, so defined, looks like a directed graph with extra structure. This is correct—and instructive. Network analysis has made important inroads into IR precisely by foregrounding relational structure over unit attributes (Hafner-Burton et al., 2009; Maoz, 2012). A network (directed graph) specifies actors and pairwise ties—who is connected to whom, and with what strength or type. But a network, by itself, says nothing about how ties *compose*. It can tell you that the United States has a trade tie with Vietnam and Vietnam has a trade tie with Laos, but it cannot, without further machinery, say what the composite relationship is. Centrality measures, clustering coefficients, and brokerage scores are all computed from the graph structure alone; they describe *patterns* of ties but do not theorize the *algebra* of how ties combine.

Moreover, as Zhukov and Stewart (2013) demonstrate, results in network-based IR research are acutely sensitive to the choice of network specification—geographic proximity, trade ties, alliance ties, and IGO co-membership can yield conflicting conclusions about the same diffusion process. In the framework developed here, each such specification corresponds to a different *category* (same objects, different morphisms), and the sensitivity is not a pathology but the expected consequence of working in different relational structures. The question of how results in one network translate to another becomes, formally, a question about *functors between categories*—a topic we take up in Section 6.

Category theory adds exactly the missing structure. Composition is built in as an axiom, and the identity and associativity conditions ensure that the resulting algebra of relations is coherent. In this sense, the move from network analysis to category theory is analogous to the move from describing a system to theorizing about it—from recording which ties exist to specifying what it means for ties to compose, and what follows from that specification.

## 2.2 Example: The trade category

**Example 2.3** (The trade category  $\mathcal{T}$ ). Let the objects of  $\mathcal{T}$  be states. For states  $A$  and  $B$ , a morphism  $f \in \text{Hom}_{\mathcal{T}}(A, B)$  represents a trade relationship

from  $A$  to  $B$ —an export flow, a trade agreement, a supply-chain dependency. There may be multiple morphisms between the same pair of states, reflecting different trade channels or agreements.

Composition captures transitive linkage: if  $f: A \rightarrow B$  represents  $A$ 's export relationship to  $B$  and  $g: B \rightarrow C$  represents  $B$ 's export relationship to  $C$ , then the composite  $g \circ f: A \rightarrow C$  represents the indirect trade relationship between  $A$  and  $C$  mediated by  $B$ . This is the formal counterpart of the intuition that global supply chains create relationships between states that do not trade directly.

The identity morphism  $\text{id}_A$  represents  $A$ 's domestic economy—the self-exchange that is always available and that leaves other trade relationships unchanged when composed with them.

To make this concrete, consider a small fragment of the trade category with four objects: the United States, China, Vietnam, and Laos. The morphisms include, among others: a massive bilateral export flow from China to the United States; an export flow from Vietnam to China driven by electronics assembly; a smaller flow from Laos to Vietnam in agricultural goods; and so on. Composition gives us the indirect trade linkages: the composite of the Laos→Vietnam and Vietnam→China morphisms is an indirect trade relationship from Laos to China mediated by Vietnam—a relationship that is real (Laotian raw materials enter Chinese supply chains via Vietnamese intermediaries) but invisible in any dataset that records only direct bilateral flows.

We will return to this four-state example throughout the paper. It is small enough to work with explicitly but rich enough to illustrate every concept we introduce: models as functors that assign data to these states and their trade links (Section 3), translations between models as natural transformations (Section 4), the relational profile of a state as a representable functor (Section 5.1), and the Yoneda lemma as a theorem about what determines a state's identity within this structure (Section 5).

It is worth pausing on what the trade category *does not* say. It specifies actors and trade relationships, but it does not yet tell us what trade *means*—whether trade produces welfare gains, shifts power balances, or generates interdependence. That interpretive step requires a *model*, which we will introduce in Section 3. The category is the structural skeleton; the model is the flesh.

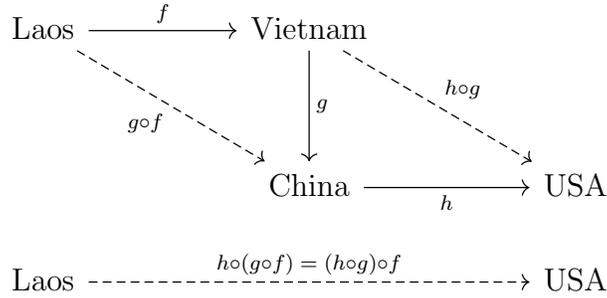


Figure 1: A fragment of the trade category  $\mathcal{T}$ . *Top*: solid arrows are direct trade morphisms; dashed arrows are composites. The composite  $g \circ f$  (Laos to China, via Vietnam) and  $h \circ g$  (Vietnam to the USA, via China) represent different intermediate mediations. *Bottom*: the fully mediated relationship from Laos to the USA. Associativity guarantees that the two ways of computing it agree:  $h \circ (g \circ f) = (h \circ g) \circ f$ .

### 2.3 Example: The alliance category

**Example 2.4** (The alliance category  $\mathcal{A}$ ). Let the objects of  $\mathcal{A}$  be states. A morphism  $f \in \text{Hom}_{\mathcal{A}}(A, B)$  represents an alliance commitment from  $A$  to  $B$ —a mutual defense pact, a security guarantee, or an informal alignment. Composition captures *extended deterrence*: if  $A$  is allied to  $B$  and  $B$  is allied to  $C$ , the composite  $g \circ f$  represents the chain of commitments linking  $A$  to  $C$ . Whether such chains are credible is an empirical and strategic question; the category records their formal existence as part of the alliance structure.

The identity morphism captures a state’s commitment to its own defense—the baseline security relationship that exists independently of any alliance.

The alliance category illustrates an important feature of the framework: the same set of states can be organized into different categories by choosing different morphism types. The United States, Japan, and Australia are the same objects in both  $\mathcal{T}$  and  $\mathcal{A}$ , but their morphisms—trade flows versus security commitments—are entirely different. This is precisely the sense in which the choice of morphisms constitutes a theoretical commitment.

### 2.4 Example: The diplomatic recognition category

**Example 2.5** (The recognition category  $\mathcal{R}$ ). Let the objects of  $\mathcal{R}$  be states and would-be states. A morphism  $f \in \text{Hom}_{\mathcal{R}}(A, B)$  represents  $A$ ’s diplomatic

recognition of  $B$  as a sovereign state.

This example is particularly revealing for several reasons. First, recognition is *non-symmetric*: the People’s Republic of China does not recognize Taiwan as a sovereign state, though Taiwan’s government maintains its own set of recognition relationships. Second, recognition is *politically constitutive* in a way that trade and alliance are not: the very status of an entity as an object in the category (a “state”) is partly determined by whether other objects recognize it. Kosovo is an object in  $\mathcal{R}$  for those states that recognize it and arguably not for those that do not. The category-theoretic framework forces this ambiguity into the open rather than hiding it behind verbal hedges.

Third, the composition axiom raises a genuinely substantive question. If  $A$  recognizes  $B$  and  $B$  recognizes  $C$ , does  $A$  bear any implied recognition relationship to  $C$ ? The answer is not obvious—and the fact that the formalism *requires* an answer is a feature, not a bug. It forces the theorist to be explicit about assumptions that are often left implicit.

## 2.5 What categories do not yet capture

A category, by itself, is an abstract relational structure. It tells us who the actors are and how they are connected, but it does not interpret those connections—it does not say what they *produce*, what they *mean*, or how they can be *measured*. The trade category  $\mathcal{T}$  records that states are linked by trade, but it is silent on whether trade brings peace, prosperity, or vulnerability.

In particular, the question “does trade cause peace?” is not a question that lives inside a single category. It is a question *between* two categories—the trade category  $\mathcal{T}$  and some conflict category—and answering it requires specifying the structural relationship between them: a functor, a shared model, or a natural transformation. The framework does not dissolve such questions; it forces them into a precise form. We return to this in Sections 3 and 6.

This is by design. A theory, in the sense intended here, is a *signature*—a specification of the basic vocabulary (actors, relations, composition) without an interpretation. To give a theory content, we need to assign concrete meaning to its abstract structure. In the language of mathematical logic, we need a *model*. In the language of category theory, we need a *functor*. That is the subject of the next section.

But before moving on, it is worth noting what the category-theoretic

framing already makes visible about existing research practice. Consider the dominant data infrastructure of quantitative IR. Datasets like the Penn World Tables or the Correlates of War National Material Capabilities index are organized as *attribute tables*: rows are states, columns are properties (GDP, military expenditure, population). Relations between states—the morphisms of the category—are either absent entirely or flattened into dyadic variables that discard composition. As Hoff and Ward (2004) observe, most empirical studies in IR assume “not only that the major actors are sovereign, but also that their relationships are portrayed in data that are modeled as independent phenomena.” The COW bilateral trade dataset, for instance, records that the United States and China have a trade relationship of a certain volume, but the mediated relationship from Laos to the United States via Vietnam and China (Figure 1) is not a record in the dataset. It must be reconstructed, if it is considered at all.

The framework developed here suggests that this is not merely a practical limitation but a theoretical one: the standard data format encodes an *attribute-first* ontology in which states are characterized by intrinsic properties and relations are secondary. A category-theoretic approach inverts this priority. Relations are the primary data; attributes are derived from them via functors. We will make this inversion precise in Section 3, where we show that familiar datasets can be understood as specific functors applied to specific categories—and that what they measure, and what they miss, depends on both choices.

### 3 What Is a Model?

A theory, in the sense of Section 2, specifies actors and relations but says nothing about what those relations produce, measure, or mean. The trade category  $\mathcal{T}$  tells us that states are connected by trade, but not whether trade enriches, empowers, or endangers them. To interpret a theory—to give it empirical or normative content—we need to assign concrete data to its abstract structure. The mathematical tool for this is the *functor*.

#### 3.1 Functors as structure-preserving maps

**Definition 3.1** (Functor). Let  $\mathcal{C}$  and  $\mathcal{D}$  be categories. A *functor*  $F: \mathcal{C} \rightarrow \mathcal{D}$  consists of:

- (i) an assignment of an object  $F(A) \in \mathcal{D}$  for each object  $A \in \mathcal{C}$ ;
- (ii) an assignment of a morphism  $F(f): F(A) \rightarrow F(B)$  in  $\mathcal{D}$  for each morphism  $f: A \rightarrow B$  in  $\mathcal{C}$ ;

subject to:

- $F(\text{id}_A) = \text{id}_{F(A)}$  for all objects  $A$ ;
- $F(g \circ f) = F(g) \circ F(f)$  for all composable morphisms  $f, g$ .

The definition says: a functor is a map from one category to another that preserves the relational structure. It sends objects to objects, morphisms to morphisms, and respects both composition and identity. Nothing is lost in translation—every relationship in the source category has an image in the target, and the way relationships compose is faithfully tracked.

## 3.2 Models as functors to Set

The general definition allows functors between any two categories, but one target category is of special importance.

**Definition 3.2** (Model of a theory). A *model* of a theory  $\mathcal{C}$  is a functor  $F: \mathcal{C} \rightarrow \mathbf{Set}$ , where  $\mathbf{Set}$  is the category whose objects are sets and whose morphisms are functions.

What does this mean concretely? A model assigns to each actor  $A$  a set  $F(A)$ —the “data” associated with  $A$  under this interpretation. These might be welfare levels, power rankings, policy positions, or observable behaviors; the choice is the model’s to make. To each relationship  $f: A \rightarrow B$ , the model assigns a function  $F(f): F(A) \rightarrow F(B)$  that specifies how the relationship transforms data—how  $A$ ’s attributes are carried, filtered, or distorted as they pass through the relation to  $B$ .

The two functorial axioms now carry substantive meaning. The condition  $F(\text{id}_A) = \text{id}_{F(A)}$  says that the self-relation leaves a state’s data unchanged—the identity morphism, whatever substantive content it carries at the level of the theory (Section 2.1), acts trivially on the model’s data. The condition  $F(g \circ f) = F(g) \circ F(f)$  says that the model respects mediated relationships:

the data transformation induced by the indirect link from  $A$  to  $C$  (via  $B$ ) is the same whether we compute it in one step or two.<sup>2</sup>

This framing inverts the usual relationship between theory and data. In standard practice, we begin with a dataset—a table of attributes—and then ask which theory best explains the patterns. Here, the theory comes first as a relational structure, and the dataset is a *consequence* of the choice of model: the functor  $F$  determines what gets measured, for whom, and how measurements at different nodes are related. Different functors applied to the same theory yield different datasets—not because the world has changed, but because the interpretive lens has.

### 3.3 Two models of the trade category

**Example 3.3** (Liberal and realist models of trade). Consider two models of the trade category  $\mathcal{T}$  from Example 2.3:

- The *liberal model*  $F: \mathcal{T} \rightarrow \mathbf{Set}$  assigns to each state  $A$  the set  $F(A)$  of possible welfare outcomes for  $A$ 's economy—real income levels, consumer surplus, gains from specialization. To each trade morphism  $f: A \rightarrow B$ , it assigns a function  $F(f): F(A) \rightarrow F(B)$  that tracks how  $A$ 's welfare position translates into welfare effects on  $B$  through the trade link. The liberal model reads trade as a mechanism of mutual enrichment.
- The *realist model*  $G: \mathcal{T} \rightarrow \mathbf{Set}$  assigns to each state  $A$  the set  $G(A)$  of relative power positions—military capacity, economic leverage, technological advantage relative to competitors. To each trade morphism  $f: A \rightarrow B$ , it assigns a function  $G(f): G(A) \rightarrow G(B)$  that maps  $A$ 's power position to the power shift induced in  $B$  by the trade relationship. The realist model reads trade as a vector of strategic dependence.

Both  $F$  and  $G$  are models of the *same* underlying theory  $\mathcal{T}$ . They agree on who the actors are and what relationships exist between them. They disagree

---

<sup>2</sup>This is Lawvere's insight. In his 1963 thesis (Lawvere, 1963), Lawvere showed that the classical notion of a “model of a theory” in mathematical logic—an interpretation that assigns sets and operations to the symbols of a formal language—is exactly a functor from the theory (understood as a category) to  $\mathbf{Set}$ . We are applying the same idea to theories of international relations. The vocabulary is different, but the logic is identical: a theory specifies structure; a model interprets that structure in the world of sets and functions.

on what those relationships *mean*—on the sets assigned to each state and the functions assigned to each trade link.

This is the precise sense in which the liberal–realist debate about trade is not a disagreement about structure but about interpretation. The two camps share a theory (states linked by trade) but operate with different models (welfare versus power). The question of whether they can be systematically compared—whether there is a coherent translation from welfare data to power data that respects the trade structure—is a question about *natural transformations*, the subject of Section 4.

### 3.4 The running example, modeled

Return to the four-state fragment of  $\mathcal{T}$ : Laos, Vietnam, China, and the United States (Figure 1). Under the liberal model  $F$ , each state is assigned a set of welfare outcomes:

$F(\text{Laos})$	welfare outcomes for Laos (agricultural export revenues, etc.)
$F(\text{Vietnam})$	welfare outcomes for Vietnam (manufacturing wages, FDI effects, etc.)
$F(\text{China})$	welfare outcomes for China (supply-chain value-added, etc.)
$F(\text{USA})$	welfare outcomes for the USA (consumer prices, import competition, etc.)

The trade morphism  $f: \text{Laos} \rightarrow \text{Vietnam}$  (agricultural exports) maps, under  $F$ , to a function  $F(f): F(\text{Laos}) \rightarrow F(\text{Vietnam})$  that tracks how Lao-tian welfare translates into Vietnamese welfare effects via the trade link. The composite morphism  $g \circ f: \text{Laos} \rightarrow \text{China}$  maps to  $F(g) \circ F(f)$ —the composition of welfare-transmission functions, capturing the indirect welfare linkage from Laos to China mediated by Vietnam.

Under the realist model  $G$ , the same states are assigned sets of power positions, and the same morphisms are assigned power-transmission functions. The *same* composite  $g \circ f$  now maps to  $G(g) \circ G(f)$ —not a chain of welfare effects but a chain of strategic dependencies. The structure is identical; the interpretation is not.

### 3.5 Strategic interaction as a multi-sorted theory

The trade, alliance, and recognition categories of Section 2 all have a single sort of object: states. But many theories in IR involve objects of fundamentally different kinds. Game-theoretic models are a natural example.

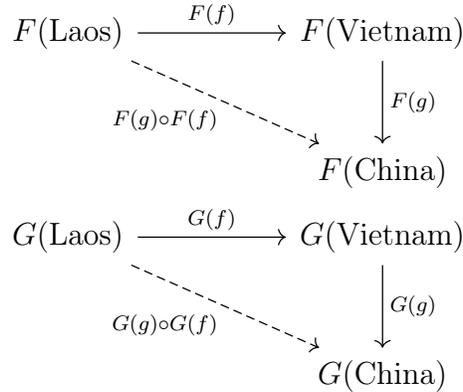


Figure 2: The liberal model  $F$  (above) and the realist model  $G$  (below) applied to the same fragment of the trade category. Both functors preserve the compositional structure: the indirect link from Laos to China decomposes identically. What differs is the content of the sets and the meaning of the functions.

**Example 3.4** (The bargaining skeleton). Consider the class of bargaining models of war. At the highest level of abstraction, these models share a common structure—a *bargaining skeleton*  $\mathcal{B}$  whose objects fall into three sorts:

- *Players*—the states engaged in the dispute;
- *The good*—the object of contention (territory, policy, status);
- *Outcomes*—the possible resolutions (a peaceful division, or war with its attendant costs and probabilities).

The morphisms encode the structure of strategic interaction: there are morphisms from players to the good (demand functions), from players to outcomes (payoff maps), and a contest morphism that determines who prevails in war.

Different bargaining models are different *functors* on this shared skeleton. Two of the cleanest examples—chosen because they share the same primitives but differ sharply in their assumptions—are the following.

**Example 3.5** (Fearon’s model as a functor). Fearon’s (1995) model sends  $\mathcal{B}$  to  $\mathbf{Set}$  as follows. The good maps to  $[0, 1]$ . Each player’s action space is a

set of offers and accept/reject decisions. The contest morphism maps to an exogenous probability  $p \in (0, 1)$ , and costs of war are exogenous parameters  $c_1, c_2 > 0$ . The payoff morphisms send player 1 to  $x$  (if peace at division  $x$ ) or to  $p - c_1$  (if war), and player 2 to  $1 - x$  or  $(1 - p) - c_2$ . Call this functor  $F$ .

**Example 3.6** (Beviá and Corchón’s model as a functor). Beviá and Corchón (2010) model the same skeleton differently. The good maps to total resources  $V = V_1 + V_2 \in \mathbb{R}_+$ . Each player’s action space now includes transfers and effort choices  $e_i$ —richer than Fearon’s accept/reject. The contest morphism maps not to an exogenous  $p$  but to an endogenous contest success function  $p_1 = \lambda e_1^\gamma / (\lambda e_1^\gamma + e_2^\gamma)$ , where  $\lambda$  measures military proficiency and  $\gamma$  measures the sensitivity of winning to effort. Costs are endogenous: war effort costs  $ke_i$ , a fixed proportion of resources committed. Call this functor  $G$ .

The two models share the same primitives—players, a contested good, a contest, payoffs—but their assumptions differ in a structured way. Fearon treats the probability of winning and the costs of war as exogenous parameters; Beviá and Corchón endogenize both through the effort choices and the contest success function. In our language:  $F$  and  $G$  are two functors on the same category  $\mathcal{B}$ , differing in what they assign to the contest and cost morphisms.

Notice the role of equilibrium. The category  $\mathcal{B}$  has a contest morphism—that is a primitive, part of the syntax of the theory. But what the contest morphism *does* is the functor’s job, and in Beviá and Corchón’s model the functor can only do its job once the game has been solved. The contest success function  $\lambda e_1^\gamma / (\lambda e_1^\gamma + e_2^\gamma)$  is a function of effort levels, not yet a concrete map between the sets that  $G$  assigns to Players and Outcomes. It becomes one only at equilibrium, when optimal effort choices  $e_1^*, e_2^*$  pin down a determinate probability  $p^*$  and determinate costs  $ke_i^*$ . Solving the model is what makes the functor well-defined: equilibrium is not a separate layer sitting on top of the model but part of the specification of  $G$ ’s assignments. Fearon’s functor  $F$  sidesteps this entirely—its assignments are exogenous parameters, so no solution concept is needed to pin them down. The difference between the two functors is therefore not just a difference in parameter values but a difference in *what kind of theoretical work is required to specify the functor at all*.

The point is general, and it clarifies a distinction that game theorists often leave implicit. The *primitives* of a model—players, a contested good,

action spaces, a payoff structure—are the objects and morphisms of a category: the syntax of the theory, specifying what the model is *about*. The *assumptions*—that the good is  $[0, 1]$ , that utility is linear, that costs are exogenous or endogenous—are the functor’s assignments: the concrete sets and functions used to *represent* that syntax in **Set**. Two models that share the same primitives but differ in their assumptions are two functors on the same category. Two models that differ in their primitives—say, one that includes an information structure and one that does not—are functors on different categories altogether, and comparing them requires a functor between the categories themselves (Section 6).

The question of whether Fearon’s model and Beviá and Corchón’s model can be systematically related—whether there is a coherent translation from one to the other—is a question about *natural transformations*, the subject of Section 4.

### 3.6 Existing datasets as implicit functors

The functor framework makes it possible to say precisely what existing IR datasets are doing—and what they are missing. The idea that a database schema is a category and a database instance is a functor to **Set** is well established in applied category theory; Fong and Spivak (2019, ch. 3) give a thorough introduction. Our contribution is to apply this lens to the data infrastructure of IR, where the implicit categorical commitments have gone unexamined.

**The Penn World Tables as a functor.** The Penn World Tables (PWT) assign to each state a vector of macroeconomic indicators: real GDP, capital stock, productivity levels, terms of trade. In our framework, this is a functor—but a very particular one. The PWT treats the international system as a *discrete category*: a collection of states with no morphisms other than identities. The functor  $F_{\text{PWT}}$  assigns to each state  $A$  the set  $F_{\text{PWT}}(A)$  of its macroeconomic data, and to each identity morphism  $\text{id}_A$  the identity function on that set. There is nothing else to assign, because the theory underlying the PWT has no relations.

This is not a criticism of the Penn World Tables—they are an extraordinary achievement in data harmonization. It is an observation about what kind of theory they implicitly encode. A dataset organized as an attribute table, with rows indexed by states and no relational structure, is a model of a

theory in which states are atoms: unconnected, self-contained, characterized entirely by intrinsic properties.

**The COW bilateral trade dataset as a functor.** The Correlates of War bilateral trade data is richer. It records, for each pair of states, a trade volume—making it a model of something closer to the trade category  $\mathcal{T}$ . But the COW data does not record composition. The trade volume between the United States and China is a direct measurement; the mediated trade relationship from Laos to the United States via Vietnam and China (the composite  $h \circ g \circ f$  in Figure 1) is not a row in the dataset. In functorial terms, the COW data is a model of the *underlying graph* of  $\mathcal{T}$ —the directed graph that records which morphisms exist—but not of the full category, because it does not track how trade relationships compose.<sup>3</sup>

**The COW National Material Capabilities index as a functor.** The NMC index (Composite Index of National Capability, or CINC score) assigns to each state a scalar summary of military expenditure, military personnel, energy consumption, iron and steel production, urban population, and total population. Like the PWT, this is a functor on a discrete category. But here the choice of model is more transparently theoretical: the CINC score is a realist model in miniature. It assigns to each state a measure of *material capability*—a quantity that makes sense only within a theory that foregrounds power. The liberal economist’s objection that GDP per capita is a better measure of national strength is not a dispute about data but about models: it is a dispute about which functor to apply to the same underlying theory.

**What these observations buy us.** Recognizing datasets as functors does two things. First, it reveals the implicit theoretical commitments embedded in data infrastructure. The shape of a dataset—attribute table versus relational table, bilateral versus multilateral, static versus compositional—is not a neutral design choice but an encoding of ontological assumptions about what kind of theory the data is meant to serve. Second, it opens the door to a precise account of what is lost when a richer relational structure is flattened

---

<sup>3</sup>More precisely, the COW bilateral trade data can be understood as a functor on the *free category* generated by the trade graph, modulo the additional assumption that only direct edges carry data. The composite morphisms exist in the free category, but the functor assigns them no independent empirical content.

into a simpler data format. The PWT functor discards all morphisms; the COW trade functor discards composition. Each discarding is a theoretical choice, and the framework developed here makes those choices explicit rather than leaving them buried in the column headers of a spreadsheet.

## 4 Translating Between Models

We now have the vocabulary to say what a theory is (a category) and what a model is (a functor to **Set**). Section 3.3 showed that the liberal and realist models of the trade category interpret the same relational structure in different ways—welfare versus power. The natural question is: can we systematically translate between them?

Not just for one state, but for *every* state, and in a way that respects the relational structure? The answer lies in the concept of a *natural transformation*—the central notion of coherent translation in category theory, and arguably the concept that motivated the entire subject.<sup>4</sup>

### 4.1 Definition and the naturality condition

**Definition 4.1** (Natural transformation). Let  $F, G: \mathcal{C} \rightarrow \mathbf{Set}$  be two models of a theory  $\mathcal{C}$ . A *natural transformation*  $\eta: F \Rightarrow G$  consists of, for each object  $A$  in  $\mathcal{C}$ , a function  $\eta_A: F(A) \rightarrow G(A)$ , such that for every morphism  $f: A \rightarrow B$  in  $\mathcal{C}$ , the following diagram commutes:

$$\begin{array}{ccc} F(A) & \xrightarrow{\eta_A} & G(A) \\ F(f) \downarrow & & \downarrow G(f) \\ F(B) & \xrightarrow{\eta_B} & G(B) \end{array}$$

That is,  $G(f) \circ \eta_A = \eta_B \circ F(f)$ .

The diagram is called the *naturality square*, and the condition it expresses is called *naturality*. In words: translating  $A$ 's data and then pushing it through  $G$ 's version of the relationship gives the same result as pushing

---

<sup>4</sup>Eilenberg and Mac Lane introduced categories and functors in 1945 primarily as scaffolding for defining natural transformations. As Mac Lane later wrote, “I didn’t invent categories to study functors; I invented them to study natural transformations” (Mac Lane, 1998, p. 29).

it through  $F$ 's version of the relationship and then translating  $B$ 's data. Translation commutes with the relational structure.

This is a strong requirement. A natural transformation is *not* merely a collection of functions  $\eta_A: F(A) \rightarrow G(A)$ , one per state, chosen independently. The naturality condition ties these functions together: they must be compatible with every morphism in the theory. Any translation that works for isolated states but breaks down when you account for their relationships is not natural—and in this framework, it is not a legitimate translation at all.

## 4.2 The naturality condition in IR

Return to the liberal model  $F$  and the realist model  $G$  of the trade category  $\mathcal{T}$  (Example 3.3). A natural transformation  $\eta: F \Rightarrow G$  would be a systematic way to convert welfare data into power data for every state, such that the conversion respects trade linkages.

Concretely, consider the trade morphism  $f: \text{Laos} \rightarrow \text{Vietnam}$  (Laotian agricultural exports to Vietnam). The naturality square for  $f$  says:

$$\begin{array}{ccc} F(\text{Laos}) & \xrightarrow{\eta_{\text{Laos}}} & G(\text{Laos}) \\ F(f) \downarrow & & \downarrow G(f) \\ F(\text{Vietnam}) & \xrightarrow{\eta_{\text{Vietnam}}} & G(\text{Vietnam}) \end{array}$$

Reading around the square:

- *Top then right:* Take Laos's welfare data. Convert it to power data (via  $\eta_{\text{Laos}}$ ). Then apply the realist trade function  $G(f)$  to see how Laos's power position affects Vietnam.
- *Left then bottom:* Take Laos's welfare data. Apply the liberal trade function  $F(f)$  to see how Laos's welfare affects Vietnam's welfare. Then convert Vietnam's welfare data to power data (via  $\eta_{\text{Vietnam}}$ ).

Naturality demands that these two paths give the same answer. Converting welfare to power *before* tracing a trade link must agree with tracing the trade link *first* and then converting. The translation cannot distort the relational structure.

This requirement has teeth. It rules out ad hoc conversions—say, a conversion that works well for large economies but fails for small ones whose

welfare positions are dominated by a single trade link. If the translation breaks for any state or any trade morphism, the naturality condition fails and no natural transformation exists between the two models.

### 4.3 When translation fails: incommensurability

The absence of a natural transformation between two models is not a deficiency of the framework; it is an informative result. It means that the two models assign meaning to the relational structure in ways that cannot be coherently reconciled.

Consider what this implies for the liberal–realist debate about trade. If no natural transformation  $F \Rightarrow G$  exists, then there is no systematic, structure-respecting way to translate welfare data into power data across the international system. The two interpretations are not merely different; they are *incommensurable* in a precise sense: the relational structure of trade imposes constraints on translation that cannot be simultaneously satisfied.

This gives formal content to a claim that has long circulated in IR metatheory: that paradigms can be genuinely incommensurable, not just different in emphasis. [Jackson \(2011\)](#) argues that different approaches to IR are grounded in distinct philosophical ontologies that resist straightforward comparison. [Monteiro and Ruby \(2009\)](#) push back, arguing that the appearance of incommensurability dissolves under sufficiently careful philosophical analysis. In the present framework, the question is neither philosophical nor verbal—it is structural. Two models of the same theory are commensurable if and only if a natural transformation exists between them. Whether one does is a mathematical question with a determinate answer, not a matter of interpretation.

*Remark 4.2.* It is worth noting that the existence of a natural transformation  $\eta: F \Rightarrow G$  does not require the translation to be invertible. A natural transformation may be one-directional: one can translate from welfare to power without being able to translate back. When  $\eta$  is invertible—when each  $\eta_A$  is a bijection—the natural transformation is called a *natural isomorphism*, and the two models are structurally indistinguishable. Natural isomorphisms are rare and informative: they say that two apparently different interpretations are, at bottom, the same. More commonly, a natural transformation exists in one direction but not the other, reflecting an asymmetry between the two models’ expressive power.

## 4.4 The running example

Let us trace the naturality condition through the full four-state fragment of Figure 1. A natural transformation  $\eta: F \Rightarrow G$  must supply four functions:

$$\begin{aligned} \eta_{\text{Laos}}: F(\text{Laos}) &\rightarrow G(\text{Laos}), & \eta_{\text{Viet}}: F(\text{Viet}) &\rightarrow G(\text{Viet}), \\ \eta_{\text{China}}: F(\text{China}) &\rightarrow G(\text{China}), & \eta_{\text{USA}}: F(\text{USA}) &\rightarrow G(\text{USA}). \end{aligned}$$

These must satisfy a naturality square for *every* trade morphism. For the three direct morphisms  $f, g, h$  and their composites, this gives:

$$\begin{aligned} G(f) \circ \eta_{\text{Laos}} &= \eta_{\text{Viet}} \circ F(f), \\ G(g) \circ \eta_{\text{Viet}} &= \eta_{\text{China}} \circ F(g), \\ G(h) \circ \eta_{\text{China}} &= \eta_{\text{USA}} \circ F(h). \end{aligned}$$

The composite morphisms  $g \circ f$  and  $h \circ g$  generate no additional constraints, because functoriality already guarantees  $G(g \circ f) = G(g) \circ G(f)$  and similarly for  $F$ . But the three equations above are already demanding: the translation functions at four states must be simultaneously compatible with three trade links.

Figure 3 displays the naturality condition as a ladder of commuting squares. Each rung connects the liberal and realist interpretations at a single state; each vertical arrow traces a trade link within one model. The ladder makes visible what “coherent translation” demands: it is not enough to convert welfare to power state by state. The conversions must interlock across every trade relationship in the system.

## 4.5 Naturality and the bargaining model

The same logic applies to the game-theoretic setting of Section 3.5, and here it yields a particularly concrete payoff. Recall the two models of the bargaining skeleton  $\mathcal{B}$ : Fearon’s functor  $F$  (Example 3.5), which treats the probability of winning  $p$  and the costs of war  $c_i$  as exogenous parameters, and Beviá and Corchón’s functor  $G$  (Example 3.6), which endogenizes both through effort choices and a contest success function.

A natural transformation  $\eta: F \Rightarrow G$  must supply a function for each object in  $\mathcal{B}$ :

- $\eta_{\text{Good}}: [0, 1] \rightarrow \mathbb{R}_+$ , mapping Fearon’s normalized pie to Beviá and Corchón’s resource space;

$$\begin{array}{ccc}
F(\text{Laos}) & \xrightarrow{\eta_{\text{L}}} & G(\text{Laos}) \\
\downarrow F(f) & & \downarrow G(f) \\
F(\text{Viet}) & \xrightarrow{\eta_{\text{V}}} & G(\text{Viet}) \\
\downarrow F(g) & & \downarrow G(g) \\
F(\text{China}) & \xrightarrow{\eta_{\text{C}}} & G(\text{China}) \\
\downarrow F(h) & & \downarrow G(h) \\
F(\text{USA}) & \xrightarrow{\eta_{\text{U}}} & G(\text{USA})
\end{array}$$

Figure 3: The naturality condition for a translation  $\eta: F \Rightarrow G$  across the four-state trade chain. Each square must commute: translating and then tracing a trade link must equal tracing the link and then translating. The entire ladder of squares must be satisfied simultaneously.

- $\eta_{\text{Player}_i}$ , translating each player’s data from one model to the other;
- $\eta_{\text{Outcome}}$ , mapping Fearon’s payoff space to Beviá and Corchón’s.

The naturality condition demands that these translations commute with the morphisms of  $\mathcal{B}$ —the demand functions, the contest morphism, and the payoff maps.

The critical square is the one involving the contest morphism. In Fearon’s model, the contest morphism maps to an exogenous probability  $p$ . In Beviá and Corchón’s, it maps to the equilibrium probability  $p^*(V_1, V_2, k, \lambda, \gamma) = \lambda e_1^{*\gamma} / (\lambda e_1^{*\gamma} + e_2^{*\gamma})$  determined by the players’ optimal effort choices. The naturality square says: translating Fearon’s player data into Beviá and Corchón’s and *then* applying the contest success function must give the same result as applying Fearon’s exogenous  $p$  and *then* translating the outcome.

$$\begin{array}{ccc}
F(\text{Players}) & \xrightarrow{\eta_{\text{Players}}} & G(\text{Players}) \\
\downarrow p & & \downarrow p^* \\
F(\text{Outcome}) & \xrightarrow{\eta_{\text{Outcome}}} & G(\text{Outcome})
\end{array}$$

This is a substantive constraint, and we can make it completely explicit. In Beviá and Corchón’s model, the first-order conditions for optimal effort are:

$$\frac{\gamma\lambda e_1^{\gamma-1} e_2^\gamma}{(\lambda e_1^\gamma + e_2^\gamma)^2} \cdot V = k \quad \text{and} \quad \frac{\gamma\lambda e_1^\gamma e_2^{\gamma-1}}{(\lambda e_1^\gamma + e_2^\gamma)^2} \cdot V = k.$$

Dividing the first by the second gives  $e_2^*/e_1^* = 1$ : at equilibrium, both players exert the same effort. It follows that  $p^* = \lambda/(\lambda + 1)$  and, crucially, that the equilibrium costs of war satisfy  $c_1^* = ke_1^* = ke_2^* = c_2^*$ . Substituting back,  $e^* = \gamma\lambda V/(k(\lambda + 1)^2)$ , so  $c^* = \gamma p^*(1 - p^*)V$ .

The constraint  $c_1^* = c_2^*$  is not an artifact of the symmetric cost parameter  $k$ . Even if we allow different marginal costs  $k_1, k_2$  for the two players, the first-order conditions give  $e_2^*/e_1^* = k_1/k_2$ , and so  $c_1^* = k_1 e_1^* = k_2 e_2^* = c_2^*$ : equilibrium war costs are equal regardless.

In Fearon’s model,  $c_1$  and  $c_2$  are free parameters—there is no reason for them to coincide. A natural transformation  $\eta: F \Rightarrow G$  can therefore exist only for Fearon triples  $(p, c_1, c_2)$  with  $c_1 = c_2$ . Every triple with  $c_1 \neq c_2$  lies outside the image of Beviá and Corchón’s equilibrium map, and no structure-preserving translation can reach it.

This tells us something precise about the relationship between the two models. Beviá and Corchón’s is *not* simply a generalization of Fearon’s with endogenous parameters; it is a model that constrains the feasible  $(p, c)$  space to a proper subset—specifically, the locus  $c_1 = c_2 = \gamma p(1 - p)V$ . The naturality condition identifies exactly where the two models are compatible and where they are not. The Fearon configurations outside that locus are the ones for which the exogenous-parameter assumption does real theoretical work—they describe scenarios that *cannot* arise from any contest game of the Beviá–Corchón type.

*Remark 4.3* (What the constraint tells a modeler). The locus  $c = \gamma p(1 - p)V$  makes visible which theoretical ingredients are needed to microfound a given Fearon configuration. Fix the simplest contest:  $\lambda = \gamma = 1$  (symmetric, proportional) and  $V = 1$ . Then  $p^* = \lambda/(\lambda + 1) = 1/2$  and  $c^* = \gamma p^*(1 - p^*)V = 1/4$ , regardless of  $k$ . (The cost parameter  $k$  scales effort— $e^* = 1/(4k)$ —but cancels out of equilibrium costs:  $c^* = ke^* = 1/4$ .) The entire image in Fearon’s parameter space is a single point:  $(p, c_1, c_2) = (1/2, 1/4, 1/4)$ .

To reach other Fearon configurations, specific theoretical commitments are required. Generating  $p \neq 1/2$  requires  $\lambda \neq 1$ : asymmetric military proficiency is the ingredient that moves the probability of winning. Generating

different cost levels requires varying  $\gamma$  (non-proportional contest technology) or  $V$  (the size of the contested good). So a modeler studying a Fearon scenario with  $p = 0.7$  and  $c_1 = c_2 = 0.1$  who wants Beviá–Corchón micro-foundations can read off the required parameters:  $\lambda = p^*/(1 - p^*) = 7/3$  and  $\gamma V = c^*/[p^*(1 - p^*)] \approx 0.476$ . The naturality condition does not merely say that a translation exists or fails; it identifies, for each point in Fearon’s parameter space, exactly what theoretical resources are needed to reach it from Beviá and Corchón’s equilibrium map—and for every point with  $c_1 \neq c_2$ , it says that no contest game of this type will suffice.

## 4.6 What we can now say

With categories, functors, and natural transformations in hand, we have the first three layers of the framework:

Category theory	IR interpretation	Section
Category $\mathcal{C}$	A theory (actors + relations)	2
Functor $F: \mathcal{C} \rightarrow \mathbf{Set}$	A model (interpretation of the theory)	3
Natural transformation $\eta: F \Rightarrow G$	A translation between models	4

This is the basic vocabulary of Lawvere’s functorial semantics ([Lawvere, 1963](#)), transplanted into international relations.

But we have not yet addressed the question that motivates this paper: what determines the identity of a state within a given theory? To answer this, we need to look at a special family of models—those that arise canonically from the theory itself, without any external interpretive choice. These are the *representable functors*, and they are the subject of the next section.

## 5 The Yoneda Lemma: Identity Is Relational

So far, every model we have considered has been an external interpretation: a functor that assigns data to the theory’s actors and relations from the outside. The liberal model reads trade as welfare; the realist model reads trade as power; the Penn World Tables read the system as a collection of disconnected attribute vectors. In each case, the interpretive choice—the functor—comes from the theorist.

But there is a family of models that arise canonically from the theory itself, without any external choice at all. These are the *representable functors*, and they encode the idea that a state’s relational profile—the totality of its relationships to every other state—is itself a model. This section introduces the relational profile, states the Yoneda lemma, and draws out its consequences. The Yoneda lemma is the central result of the paper, and its content is this: every datum that any model assigns to a state is not a local fact about that state but a global structural commitment—a complete, coherent reading of that state’s position in the system.

## 5.1 The relational profile

**Definition 5.1** (Representable functor). For any object  $X$  in a category  $\mathcal{C}$ , the *representable functor*  $\text{Hom}_{\mathcal{C}}(-, X): \mathcal{C}^{\text{op}} \rightarrow \mathbf{Set}$  is defined by:

- (i) On objects:  $\text{Hom}(-, X)$  sends each object  $A$  to the set  $\text{Hom}(A, X)$  of all morphisms from  $A$  to  $X$ .
- (ii) On morphisms: given  $g: B \rightarrow A$  in  $\mathcal{C}$ , the function  $\text{Hom}(g, X): \text{Hom}(A, X) \rightarrow \text{Hom}(B, X)$  sends  $f \mapsto f \circ g$  (precomposition with  $g$ ).

What does this say? For a fixed state  $X$ , the functor  $\text{Hom}(-, X)$  assigns to every other state  $A$  the set of all relationships that  $A$  bears to  $X$ . It collects, for the entire system, the answer to the question: “how does the world relate to  $X$ ?”

This is  $X$ ’s *relational profile*—not an intrinsic property of  $X$ , but the complete record of how every other actor in the system is connected to it.<sup>5</sup>

The relational profile is itself a functor to **Set**—a model in the sense of Definition 3.2. But it is a model of a very particular kind. Unlike the liberal or realist models, it makes no interpretive choice about what trade *means*. It simply records, for every state in the system, the full inventory of relationships pointing toward  $X$ . It is the most theoretically innocent model available: the model that says “to understand  $X$ , look at how the entire world relates to it.”

---

<sup>5</sup>The notation  $\mathcal{C}^{\text{op}}$  indicates the *opposite category*, in which all morphisms are formally reversed. This technical detail ensures that the Hom-functor is covariant (a morphism  $g: B \rightarrow A$  in  $\mathcal{C}$  induces a function  $\text{Hom}(A, X) \rightarrow \text{Hom}(B, X)$  going in the “right” direction). The reader who finds this confusing may safely ignore it; the intuition— $\text{Hom}(-, X)$  collects all relationships pointing toward  $X$ —is what matters.

**The profile as a working functor.** The definition says that  $\text{Hom}(-, X)$  acts not only on objects but on morphisms, and this is where the profile acquires its structural force. Consider the trade morphism  $g: \text{Vietnam} \rightarrow \text{China}$  in  $\mathcal{T}$ . The profile  $\text{Hom}(-, \text{USA})$  sends  $g$  to the function

$$\text{Hom}(g, \text{USA}): \text{Hom}(\text{China}, \text{USA}) \rightarrow \text{Hom}(\text{Vietnam}, \text{USA})$$

that takes each trade channel from China to the USA and extends it backward through  $g$ : the direct channel  $h: \text{China} \rightarrow \text{USA}$  maps to the composite  $h \circ g: \text{Vietnam} \rightarrow \text{USA}$ . The profile registers that any trade relationship arriving at the USA from China can be pulled back, via the Vietnam–China link, to a trade relationship arriving from Vietnam. This is precomposition, and it is what makes the profile a functor rather than a mere list: it tracks not just which relationships exist but how they propagate through the network.

## 5.2 Two profiles in the running example

**Example 5.2** (The trade profile of the United States). In the four-state fragment of  $\mathcal{T}$  (Figure 1), the representable functor  $\text{Hom}_{\mathcal{T}}(-, \text{USA})$  assigns:

$$\begin{aligned} \text{Hom}(\text{Laos}, \text{USA}) &= \{h \circ g \circ f\}, \\ \text{Hom}(\text{Viet}, \text{USA}) &= \{h \circ g\}, \\ \text{Hom}(\text{China}, \text{USA}) &= \{h\}, \\ \text{Hom}(\text{USA}, \text{USA}) &= \{\text{id}_{\text{USA}}\}. \end{aligned}$$

Every state in the fragment has a nonempty set of trade channels pointing toward the USA. The profile is full: the USA is an actor to whom the entire network converges.

**Example 5.3** (The trade profile of Laos). The profile  $\text{Hom}_{\mathcal{T}}(-, \text{Laos})$  tells a different story:

$$\begin{aligned} \text{Hom}(\text{Laos}, \text{Laos}) &= \{\text{id}_{\text{Laos}}\}, \\ \text{Hom}(\text{Viet}, \text{Laos}) &= \emptyset, \\ \text{Hom}(\text{China}, \text{Laos}) &= \emptyset, \\ \text{Hom}(\text{USA}, \text{Laos}) &= \emptyset. \end{aligned}$$

No other state in the fragment has a trade channel pointing toward Laos. The only relationship Laos bears to itself is the identity—the domestic economy. The profile is nearly empty: Laos is a source in the trade network, not a sink.

The asymmetry between these two profiles is structural, not accidental. The USA is an actor to whom many relationships converge; Laos is an actor from whom relationships emanate but to whom few return. The profiles are different because the states *are* different—not in their intrinsic properties (the category says nothing about intrinsic properties) but in how they sit within the relational structure. And the difference is visible at a glance:  $\text{Hom}(\text{China}, \text{USA}) = \{h\}$  while  $\text{Hom}(\text{China}, \text{Laos}) = \emptyset$ .

**Example 5.4** (Relational profiles in the alliance category). In the alliance category  $\mathcal{A}$ , the profile  $\text{Hom}_{\mathcal{A}}(-, \text{USA})$  collects all alliance commitments pointing toward the United States: NATO members’ commitments, bilateral defense treaties with Japan, South Korea, Australia, the Philippines, and so on. The profile  $\text{Hom}_{\mathcal{A}}(-, \text{Switzerland})$  is nearly empty—Switzerland’s defining foreign-policy feature is the paucity of alliance morphisms converging on it. In the alliance category, Switzerland and the United States are distinguished precisely by the difference in their relational profiles.

### 5.3 Reading a profile through a model

We now have two kinds of model on the table, and the contrast between them is the key to everything that follows.

On one side is the relational profile  $\text{Hom}(-, X)$ . This is the maximally innocent model: it records every relationship that every state bears to  $X$ , and nothing else. It makes no claim about what trade means, whether it enriches or endangers, whether it shifts power or distributes welfare. It is the theory talking about itself—pure relational structure, with no interpretive overlay.

On the other side is a model like  $F$ —the liberal model, the realist model, or any other functor to **Set**. This is how the theorist *wanted* to talk about international relations. It carries an interpretation: trade means welfare, or trade means power, or trade means something else entirely. It assigns concrete data to states and concrete functions to relationships.

A natural transformation  $\eta: \text{Hom}(-, X) \Rightarrow F$  is the act of *reading* the first through the second. It takes the raw relational record—all the relations, with no spin—and feeds it into the model the theorist cares about. For each state  $A$ , the component  $\eta_A: \text{Hom}(A, X) \rightarrow F(A)$  takes every relationship that  $A$  bears to  $X$  and translates it into a datum in  $F$ ’s language: a welfare level, a power position, a data point. And it does this coherently: if a

relationship  $A \rightarrow X$  factors through another state  $B$ , the translation at  $A$  must be compatible with the translation at  $B$ , because the naturality condition ties them together (Section 4.1).

So  $\eta$  is the act of interpretation itself—the moment where raw relational structure becomes substantive content. It answers the question: “given the way the entire world relates to  $X$ , what does model  $F$  have to say about each state?”

This framing makes the Yoneda lemma’s content legible before we state it. The lemma says that every such act of interpretation—every coherent way of reading  $X$ ’s relational profile through  $F$ —is completely determined by a single element: what  $F$  assigns to  $X$  itself. One datum at  $X$ , and the entire reading is fixed. There is no room for further choice, because the relational structure propagates consequences to every other state in the system.

## 5.4 The Yoneda lemma

**Theorem 5.5** (The Yoneda Lemma). *Let  $\mathcal{C}$  be a category,  $F: \mathcal{C}^{\text{op}} \rightarrow \mathbf{Set}$  a functor, and  $X$  an object of  $\mathcal{C}$ . There is a bijection, natural in both  $X$  and  $F$ :*

$$\text{Nat}(\text{Hom}_{\mathcal{C}}(-, X), F) \cong F(X).$$

In words: the set of all natural transformations from  $X$ ’s relational profile into any model  $F$  is in exact, one-to-one correspondence with what  $F$  assigns to  $X$ .

*Proof sketch.* Given a natural transformation  $\eta: \text{Hom}(-, X) \Rightarrow F$ , evaluate its component at  $X$  on the identity morphism:  $\eta_X(\text{id}_X) \in F(X)$ . This defines a map  $\text{Nat}(\text{Hom}(-, X), F) \rightarrow F(X)$ . Conversely, given an element  $a \in F(X)$ , define  $\eta_A(f) = F(f)(a)$  for each  $f \in \text{Hom}(A, X)$ . Naturality of  $\eta$  follows from functoriality of  $F$ , and the two constructions are inverse to each other.  $\square$

The proof is short—a few lines of diagram chasing—but the content is deep. The identity morphism  $\text{id}_X$  acts as a “universal probe”: every natural transformation out of  $\text{Hom}(-, X)$  is completely determined by what it does to  $\text{id}_X$ . This is the formal expression of the idea that  $X$ ’s self-relation, combined with the relational structure, pins down everything.

**What the bijection means.** The bijection has two directions, and both matter.

The direction  $\text{Nat}(\text{Hom}(-, X), F) \rightarrow F(X)$  says: every coherent reading of  $X$ 's relational profile through model  $F$  collapses to a single element of  $F(X)$ . You cannot distribute consequences across the system without committing to a specific datum about  $X$  itself. Every interpretation has a seed.

The direction  $F(X) \rightarrow \text{Nat}(\text{Hom}(-, X), F)$  says the converse, and this is the direction with teeth. Every element  $a \in F(X)$  is not a local fact about  $X$ . It is a *global structural commitment*—a complete, coherent reading of the entire system, compressed into a single datum.

Consider what this means concretely. To say “the USA has welfare level  $a$ ” sounds like a statement about the USA. It is not. It is a statement about the *entire trade network as seen from the USA's position*, because the rule  $\eta_A^a(f) = F(f)(a)$  propagates  $a$  outward through every trade channel converging on the USA. China's welfare assessment is not a separate parameter to be estimated independently; it is  $F(h)(a)$ , forced by  $a$  and the liberal model's transmission function for the direct link  $h$ . Vietnam's is  $F(g)(F(h)(a))$ , forced by the two-step chain. Laos's is  $F(f)(F(g)(F(h)(a)))$ , forced by the full three-step chain. Each of these is a *consequence* of the single datum  $a$ , propagated through the relational structure by the model's own rules. There are no free parameters left.

And the bijection says that this propagation is exhaustive and faithful. Every coherent reading of  $X$ 's relational profile through  $F$  is generated by some element of  $F(X)$ , and no two elements generate the same reading. The local datum and the global interpretation are the same object, viewed from different ends: the element  $a \in F(X)$  is the natural transformation  $\eta^a$ , and the natural transformation  $\eta^a$  is the element  $a$ . This is not a metaphor. It is the content of the theorem.

## 5.5 The Yoneda lemma at work

Let us see the bijection concretely in the four-state trade category, taking  $X = \text{USA}$ .

**The liberal model.** Take  $F$  to be the liberal model of Example 3.3. A natural transformation  $\eta: \text{Hom}(-, \text{USA}) \Rightarrow F$  must assign, for each state  $A$ ,

a function  $\eta_A: \text{Hom}(A, \text{USA}) \rightarrow F(A)$  that is compatible with trade morphisms. The Yoneda lemma says all such transformations are determined by a single element:  $\eta_{\text{USA}}(\text{id}_{\text{USA}}) \in F(\text{USA})$ .

Pick any welfare outcome  $a \in F(\text{USA})$ —say, a particular level of consumer welfare. This determines a unique natural transformation  $\eta^a$  by the rule:

$$\eta_A^a(f) = F(f)(a) \quad \text{for each } f \in \text{Hom}(A, \text{USA}).$$

Tracing the consequences:

- $\eta_{\text{China}}^a(h) = F(h)(a)$ : the welfare effect on China of its direct trade link to the USA, starting from the USA’s welfare level  $a$ .
- $\eta_{\text{Viet}}^a(h \circ g) = F(g)(F(h)(a))$ : the welfare effect on Vietnam, obtained by first tracing  $a$  through China and then through Vietnam.
- $\eta_{\text{Laos}}^a(h \circ g \circ f) = F(f)(F(g)(F(h)(a)))$ : the welfare effect on Laos, at the end of the fully mediated chain.

A single welfare datum for the USA generates a complete welfare distribution across the system. A different datum  $a' \neq a$  generates a different distribution—different at every state where the transmission functions  $F(h)$ ,  $F(g)$ ,  $F(f)$  are not constant. The bijection  $\text{Nat}(\text{Hom}(-, \text{USA}), F) \cong F(\text{USA})$  is now visible: the USA’s welfare level is not an independent assignment but a compressed summary of the entire system’s welfare structure, as seen from the USA’s position.

**The realist model.** Now take  $G$  to be the realist model. The same relational profile  $\text{Hom}(-, \text{USA})$  is the source, but the target has changed. The Yoneda lemma gives  $\text{Nat}(\text{Hom}(-, \text{USA}), G) \cong G(\text{USA})$ : each power position  $b \in G(\text{USA})$  generates a unique natural transformation  $\eta^b$  by the rule  $\eta_A^b(f) = G(f)(b)$ . Now  $\eta_{\text{China}}^b(h) = G(h)(b)$ : the power shift induced in China by its trade link to the USA, starting from the USA’s power position  $b$ . Where the liberal model distributes welfare, the realist model distributes power—from the same relational profile, through different transmission functions.

**Why the profile of Laos works differently.** The Yoneda lemma applies equally to Laos, but the bijection looks very different. Since  $\text{Hom}(\text{Viet}, \text{Laos}) =$

$\text{Hom}(\text{China}, \text{Laos}) = \text{Hom}(\text{USA}, \text{Laos}) = \emptyset$ , a natural transformation  $\eta: \text{Hom}(-, \text{Laos}) \Rightarrow F$  has nothing to assign at Vietnam, China, or the USA—the empty set admits only the empty function. The entire natural transformation is determined by  $\eta_{\text{Laos}}(\text{id}_{\text{Laos}}) \in F(\text{Laos})$ , but it carries no information to any other state. The liberal model’s welfare assessment of Laos does not propagate through the system, because no trade channels converge on Laos through which consequences could propagate.

This is the formal expression of a structural fact: in this fragment of the trade network, Laos is a peripheral actor whose relational profile gives models little to work with. The USA’s welfare level is a global commitment; Laos’s welfare level, in this fragment, is a local island. The difference is not in the states’ intrinsic properties but in the richness of their relational profiles—and the Yoneda lemma makes the difference precise.

**What datasets lose.** The Laos example is a local thinning of the relational profile: one state has few morphisms converging on it. But the datasets discussed in Section 3.6 represent a global thinning—a flattening of the entire theory.

Consider the Penn World Tables. The PWT treats the international system as a discrete category: states with no morphisms other than identities. In a discrete category, the relational profile  $\text{Hom}(-, X)$  is trivial:  $\text{Hom}(X, X) = \{\text{id}_X\}$ , and  $\text{Hom}(A, X) = \emptyset$  for all  $A \neq X$ . The Yoneda bijection  $\text{Nat}(\text{Hom}(-, X), F) \cong F(X)$  still holds, but it says almost nothing: the natural transformation picks out a single element of  $F(X)$  and propagates it nowhere, because there are no morphisms to propagate it through. Every state’s datum is a free parameter, unconstrained by any other state’s datum. “The USA has GDP  $a$ ” *really is* a local fact in the PWT’s implicit theory, because the theory contains no relational structure to make it anything more.

Now embed the same states in the trade category  $\mathcal{T}$ . The same datum  $a$  is no longer free; it is a global structural commitment, propagating through every trade channel to every other state. The difference between the two readings of “the USA has GDP  $a$ ” is not in the data but in the theory. The PWT’s attribute table treats the number as a property of the USA. The trade category, read through the Yoneda lemma, treats it as a compressed description of the USA’s entire position in the world economy.

The COW bilateral trade data sits between these extremes. It records direct trade flows—morphisms—but not composition: the mediated relation-

ship from Laos to the USA via Vietnam and China is not a row in the dataset. In the language of Section 3.6, the COW data is a functor on the underlying graph of  $\mathcal{T}$ , not the full category. The Yoneda lemma applies to this weaker structure, but its bite is reduced: a datum propagates along direct links but not along chains. What the COW data discards is precisely the compositional structure that allows the Yoneda lemma to propagate consequences through the full network.

Each of these is a choice about how much relational structure to encode in the theory, and the Yoneda lemma makes the cost of each choice legible. The more morphisms the theory contains, the more each datum is constrained by every other. The fewer morphisms, the more each datum floats free. Attribute tables are the limiting case: maximum freedom, zero relational constraint, and a Yoneda lemma with nothing to say.

## 5.6 The Yoneda embedding

The Yoneda lemma has an immediate and powerful corollary.

**Corollary 5.6** (The Yoneda Embedding). *The functor  $\mathbf{y}: \mathcal{C} \rightarrow [\mathcal{C}^{\text{op}}, \mathbf{Set}]$  defined by  $X \mapsto \text{Hom}(-, X)$  is fully faithful. That is, for any objects  $X$  and  $Y$  in  $\mathcal{C}$ :*

$$\text{Hom}_{\mathcal{C}}(X, Y) \cong \text{Nat}(\text{Hom}(-, X), \text{Hom}(-, Y)).$$

*In particular,  $X \cong Y$  if and only if  $\text{Hom}(-, X) \cong \text{Hom}(-, Y)$ .*

This is the punchline. Two states are isomorphic—structurally identical within the theory—if and only if they have identical relational profiles. There is no identity beyond relations.

**The embedding in the running example.** The profiles computed in Examples 5.2 and 5.3 make this concrete. The USA and Laos cannot be isomorphic in  $\mathcal{T}$ , because their profiles disagree:  $\text{Hom}(\text{China}, \text{USA}) = \{h\}$  while  $\text{Hom}(\text{China}, \text{Laos}) = \emptyset$ . A natural isomorphism  $\text{Hom}(-, \text{USA}) \cong \text{Hom}(-, \text{Laos})$  would require, at the component for China, a bijection between a one-element set and the empty set—which does not exist. The Yoneda embedding thus distinguishes the USA from Laos on purely relational grounds: not because they differ in GDP or population or regime type, but because the morphisms converging on them are different.

More subtly, the Yoneda embedding says that the *only* way two states can be isomorphic is if their profiles agree at *every* state in the system. Even a single discrepancy—a single state  $A$  for which  $\text{Hom}(A, X)$  and  $\text{Hom}(A, Y)$  differ—is enough to distinguish them. Identity is overdetermined by relational data: there are as many independent witnesses to the difference between two non-isomorphic states as there are states at which their profiles disagree.

## 5.7 Philosophical implications

**Against essentialism.** The Yoneda embedding is a theorem against essentialism. It says that there are no intrinsic properties of a state—no hidden essence—that are not already captured by its relational profile. If two states are relationally identical (in a given theory), they are identical, full stop. There is no room for a residual “nature” or “character” that might distinguish them beyond their relations.

**Identity is model-dependent.** There is no bare “identity of  $A$ ”—only what the liberal model assigns to  $A$ , what the realist model assigns to  $A$ , and so on.  $F(A) \neq G(A)$  in general: the welfare profile of the United States is not the same object as its power profile. This much is obvious from the definition of a functor. What is not obvious—and what the Yoneda lemma makes precise—is that every such model-specific identity is extracted from the *same* relational profile  $\text{Hom}(-, A)$  via a natural transformation. The liberal model’s welfare assessment  $F(A)$  is obtained by mapping  $A$ ’s relational profile into  $F$ ; the realist model’s power assessment  $G(A)$  is obtained by mapping the same profile into  $G$ . Different models yield different identities, but none of them is arbitrary: each is a projection of  $A$ ’s relational position through an interpretive lens that must respect the relational structure at every point. Identity is not only relational but perspectival—and the perspectives are disciplined by naturality.

**For constructivism.** Wendt’s central claim—that “identities are constituted by interaction” (Wendt, 1992)—is now a theorem, not a philosophical stance. The Yoneda embedding proves that within any relational theory, the identity of an actor is constituted by its relationships to all other actors. This does not settle the broader philosophical debate (the Yoneda lemma

is a mathematical result, not a metaphysical one), but it shows that the constructivist claim is the *correct structural consequence* of taking relations seriously. Any framework that foregrounds relational structure and satisfies the minimal axioms of a category will, as a matter of mathematical necessity, yield the conclusion that identity is relational.

**Against the like-units assumption.** Waltz’s structural realism treats states as “like units”—functionally undifferentiated actors whose differences lie only in capabilities (Waltz, 1979). The Yoneda embedding says the opposite: states are distinguished *precisely* by their relational profiles. The United States and Laos are not like units that happen to differ in GDP; they are fundamentally different objects in the trade category because the morphisms converging on them are different. Waltz was halfway right—structure matters—but wrong to conclude that units are interchangeable within it. The structure itself differentiates them.

**States as canonical points in the space of models.** The Yoneda embedding  $\mathbf{y}: \mathcal{C} \hookrightarrow [\mathcal{C}^{\text{op}}, \mathbf{Set}]$  maps each state to its relational profile, which is a functor—a model. This means that the states of the theory are themselves models: canonical, distinguished points in the space of all models. A state is not just an object that models describe; it is, in a precise sense, a *universal model of itself*.

The space of all models  $[\mathcal{C}^{\text{op}}, \mathbf{Set}]$  is vast: it contains the liberal model, the realist model, the PWT functor, and every other conceivable interpretation of the theory. The Yoneda embedding picks out, within this space, the models that correspond to actual states—the representable functors. Every other model’s assessment of  $X$  is obtained from the representable functor  $\text{Hom}(-, X)$  via a natural transformation (this is exactly what the Yoneda bijection says). So  $\text{Hom}(-, X)$  is the most complete model of  $X$  available within the theory: it contains everything that any model could extract, and every model extracts from it by projection.

For international relations, this inverts a familiar picture. We are accustomed to thinking that models describe states—that the liberal model tells us something about the USA, and the realist model tells us something else, and the USA sits there passively, waiting to be described. The Yoneda embedding says this is backwards. The USA—understood as the representable functor  $\text{Hom}(-, \text{USA})$ —is itself a model, and it is the richest model of itself

that the theory can produce. The liberal model does not describe the USA from the outside; it *projects* the USA’s own relational self-description into a particular vocabulary. The realist model projects the same self-description into a different vocabulary. The state is not the described; it is the source from which all descriptions are derived.

This is why the section is titled as it is. Identity is relational not as a philosophical preference but as a structural fact: in any theory that satisfies the axioms of a category, the Yoneda embedding proves that a state’s identity *is* its relational profile, that every model’s assessment of the state is a projection of that profile, and that no two distinct profiles can be confused by any model whatsoever. The mathematics leaves no room for an alternative.

## 6 Comparing Theories

The previous sections developed the framework within a single theory: a category  $\mathcal{C}$ , its models (functors to **Set**), translations between models (natural transformations), and the Yoneda lemma that pins identity to relational profile. But international relations does not consist of a single theory. The trade category  $\mathcal{T}$ , the alliance category  $\mathcal{A}$ , and the recognition category  $\mathcal{R}$  are three different theories of the same system of states. They share the same actors but differ in the relations they foreground. How do we compare them?

### 6.1 Functors between theories

The same tool that gives us models—the functor—also gives us theory translations.

**Definition 6.1** (Theory translation). A *theory translation* from  $\mathcal{C}$  to  $\mathcal{D}$  is a functor  $P: \mathcal{C} \rightarrow \mathcal{D}$ : an assignment of objects to objects and morphisms to morphisms that preserves composition and identity.

The definition is identical to Definition 3.1, but the conceptual shift is important. When the target is **Set**, a functor interprets a theory in the world of data. When the target is another theory, a functor translates one relational vocabulary into another.

**Example 6.2** (From trade to alliance). Suppose that major trade partnerships tend to generate security commitments—that states which trade

heavily develop overlapping interests that eventually crystallize into alliance obligations. This hypothesis can be expressed as a functor  $P: \mathcal{T} \rightarrow \mathcal{A}$  that maps each state to itself (the actors are the same) and each trade morphism to the alliance commitment it is hypothesized to induce. The functorial axioms impose constraints on this hypothesis: if the trade link from Laos to Vietnam and the link from Vietnam to China together generate a composite alliance commitment from Laos to China, then  $P$  must map the composite trade morphism  $g \circ f$  to the composite alliance morphism  $P(g) \circ P(f)$ —the induced commitment must respect the chain of relationships, not just the bilateral links.

Not every hypothesized connection between theories will satisfy these constraints. If trade links generate alliance commitments bilaterally but the commitments do not compose transitively, then no functor  $P: \mathcal{T} \rightarrow \mathcal{A}$  exists that captures the relationship. The failure is informative: it says that the two theories’ relational structures are incompatible in a specific way. The trade category composes transitively by assumption; if the hypothesized alliance commitments do not, the translation breaks.

## 6.2 Pulling back models

Theory translations do more than relate relational structures. They transport models.

If  $P: \mathcal{C} \rightarrow \mathcal{D}$  is a functor between theories and  $F: \mathcal{D} \rightarrow \mathbf{Set}$  is a model of  $\mathcal{D}$ , then the composite  $F \circ P: \mathcal{C} \rightarrow \mathbf{Set}$  is a model of  $\mathcal{C}$ . It assigns to each object of  $\mathcal{C}$  whatever  $F$  assigns to its image in  $\mathcal{D}$ , and to each morphism the corresponding function. The composite is called the *pullback* of  $F$  along  $P$ .

This is how cross-theory interpretation works. A conflict model  $F: \mathcal{D} \rightarrow \mathbf{Set}$  assigns conflict probabilities to states and conflict-transmission functions to the relations in a conflict category  $\mathcal{D}$ . If a functor  $P: \mathcal{T} \rightarrow \mathcal{D}$  translates trade relations into conflict relations, then  $F \circ P$  is a model of the trade category that reads trade as conflict: it assigns to each state the conflict probability induced by its trade position, and to each trade morphism the corresponding conflict-transmission function.

### 6.3 Does trade cause peace?

Section 2.5 observed that the question “does trade cause peace?” does not live inside a single category. We can now say precisely where it lives.

The question requires three ingredients:

- (i) A trade category  $\mathcal{T}$  and a conflict category  $\mathcal{D}$  (two theories of the same system of states).
- (ii) A functor  $P: \mathcal{T} \rightarrow \mathcal{D}$  that translates trade relations into conflict relations (a specific causal hypothesis).
- (iii) A model  $F: \mathcal{D} \rightarrow \mathbf{Set}$  that interprets conflict relations as data—probability of armed dispute, crisis escalation, or whatever the analyst measures.

The composite  $F \circ P: \mathcal{T} \rightarrow \mathbf{Set}$  is then a model of trade that reads trade as conflict. Different functors  $P$  encode different causal stories. The liberal hypothesis that trade promotes peace corresponds to one  $P$ —one that maps trade morphisms to conflict-dampening relationships. The realist hypothesis that trade dependence generates vulnerability corresponds to a different  $P$ —one that maps trade morphisms to conflict-exacerbating relationships.

The two hypotheses are not merely different verbal accounts. They are different functors, and they generate different pullback models  $F \circ P_{\text{lib}}$  and  $F \circ P_{\text{real}}$  of the trade category. Whether a natural transformation exists between these pullback models—whether the liberal and realist readings of trade-as-conflict are coherently translatable—is a determinate structural question, answerable by the methods of Section 4.

The framework does not answer the empirical question of which functor  $P$  is correct. It clarifies what kind of object the question is about: a functor between theories, subject to the compositional constraints of the source and target categories.

### 6.4 Reducibility and embedding

The properties of a theory translation  $P: \mathcal{C} \rightarrow \mathcal{D}$  characterize the structural relationship between the two theories.

- $P$  is *faithful* if it preserves distinctions: whenever  $\mathcal{C}$  has two different morphisms  $f \neq g$  between the same pair of objects,  $P(f) \neq P(g)$  in

$\mathcal{D}$ . A faithful translation says that the target theory can express every relational distinction that the source makes.

- $P$  is *full* if the source theory captures every relation in the target: for any morphism  $h: P(A) \rightarrow P(B)$  in  $\mathcal{D}$ , there exists a morphism  $f: A \rightarrow B$  in  $\mathcal{C}$  with  $P(f) = h$ . A full translation says that the source theory is relationally complete with respect to the target.
- $P$  is *fully faithful* if both conditions hold. This is an *embedding*: the source theory is a sub-theory of the target, with no relational information lost or gained.

These distinctions bear on long-standing debates about theory reduction in IR. Can the alliance category be reduced to the trade category—do all alliance commitments arise from trade relationships? This is the claim that a full functor  $P: \mathcal{T} \rightarrow \mathcal{A}$  exists. Can trade relationships be distinguished within the alliance framework—do alliance commitments preserve the fine-grained structure of trade? This is the claim that  $P$  is faithful.

The answers are empirical, but the framework makes the structural content of the claims precise. A claim of reducibility is a claim about a specific functor with specific properties, not a vague assertion that “everything is really about trade” or “alliances subsume trade.”

## 6.5 What this section leaves open

The tools introduced here—functors between theories, pullback models, faithful and full translations—handle the simplest form of inter-theory comparison: theories that share the same actors and differ in their relations. More complex comparisons require additional structure.

When two theories do not share the same actors—say, a state-level theory and a theory that includes non-state actors, or a complete-information bargaining model and an incomplete-information extension (Section 3.5)—the relationship between them involves functors that are not the identity on objects. The primitives differ, and the comparison must map not just relations but actors.

When the question is not translation but complementarity—whether two theories illuminate different aspects of the same phenomenon in a way that is more than additive—the relevant structure is an *adjunction*: a pair of functors  $P: \mathcal{C} \rightarrow \mathcal{D}$  and  $Q: \mathcal{D} \rightarrow \mathcal{C}$  satisfying a universal property that

formalizes the idea of “best approximation in one theory of a construct in the other.” Adjunctions are ubiquitous in mathematics and have been applied to knowledge representation (Spivak, 2014), but their use in IR metatheory remains unexplored.

These extensions point toward a richer metatheoretic toolkit. The present paper has focused on what can be said within a single theory, using the Yoneda lemma. The systematic comparison of theories is a natural next step, and the categorical framework is designed to support it.

## 7 Discussion

The preceding sections have built a framework in which theories are categories, models are functors, translations between models are natural transformations, and the Yoneda lemma proves that identity is relational. This section takes stock: what does the framework buy, what does it not buy, and where might it go?

### 7.1 Relation to existing IR metatheory

**Wendt and constructivism.** The paper’s central result—that identity is constituted by relational profile—vindicates Wendt’s core claim in a precise sense. “Anarchy is what states make of it” (Wendt, 1992) can now be read as a statement about models: the structure of the international system (the category) underdetermines its meaning, which is supplied by the choice of functor. “Identities are constituted by interaction” (Wendt, 1999) is the Yoneda embedding applied to social categories.

But the framework also clarifies where Wendt’s argument needed sharpening. Wendt’s constructivism is often criticized for vagueness about what “constitution” means and how it differs from causal influence (Jackson, 2011). The Yoneda lemma gives “constitution” a precise content:  $X$ ’s identity is not *caused* by its relational profile; it *is* its relational profile, in the sense that the two are in natural bijection. The distinction between constitutive and causal claims, which has generated extensive philosophical debate in IR (Wendt, 1998), is here a distinction between the Yoneda embedding (constitutive: identity is profile) and a functor between categories (causal: one relational structure produces another, as in Section 6.3).

**Waltz and structural realism.** Waltz’s *Theory of International Politics* (Waltz, 1979) makes two moves that the framework engages directly. First, Waltz insists that the structure of the international system matters—that the arrangement of units, not just their properties, shapes outcomes. This is correct, and it is precisely the claim that the relational structure (the category) constrains what models can say (the Yoneda lemma).

Second, Waltz treats states as “like units”—functionally undifferentiated actors distinguished only by capabilities. This is where the Yoneda embedding pushes back. States are not like units; they are distinguished by their relational profiles, which encode far more than material capabilities. The United States and Laos differ not because one has more GDP but because the morphisms converging on them are structurally different (Section 5.5). Waltz’s framework captures the importance of structure but discards the information that structure carries about the differentiation of units.

**Jackson and metatheory.** Jackson (2011) argues that IR scholars operate within distinct “philosophical ontologies” that shape what counts as knowledge. The framework developed here does not resolve Jackson’s pluralism, but it gives it a formal counterpart. Different theories (categories) reflect different ontological commitments about which relations are primitive. Different models (functors) on the same theory reflect different interpretive commitments about what those relations mean. The question of commensurability between paradigms—which Jackson treats as a philosophical question about the compatibility of ontologies—becomes a mathematical question about the existence of natural transformations (Section 4.3) or functors between categories (Section 6).

This does not make the philosophical questions disappear. The choice of which category to work in—whether to foreground trade, alliance, recognition, or some other relational type—remains a substantive theoretical commitment that the mathematics cannot adjudicate. But once that choice is made, the framework constrains what can coherently be said within it, and the Yoneda lemma is the tightest of those constraints.

**Kurki and causal pluralism.** Kurki (2008) argues that IR’s dominant notion of causation—Humean regularity—is too narrow, and that the discipline should embrace a pluralism of causal concepts including constitutive, dispositional, and structural causation. The framework offers a natural

home for this pluralism. Humean causation, formalized as statistical regularities between variables, corresponds to models on discrete or near-discrete categories—attribute-table data with minimal relational structure. Structural and constitutive causation correspond to models on richer categories, where the Yoneda lemma’s propagation of consequences through relational structure provides a precise mechanism for how position constitutes identity. The “datasets lose” analysis of Section 5.5 makes visible exactly what is discarded when a relationally rich theory is compressed into a Humean data format.

## 7.2 Limitations

The framework has clear boundaries, and honesty about them is important.

**Structure, not dynamics.** Categories, as presented here, are static: they encode a fixed set of actors and relations. The international system is not static. States enter and exit, relations form and dissolve, and the relational structure itself changes over time. The framework captures the structure of the system at a given moment (or under a given theoretical idealization) but does not model the processes by which that structure evolves.

**No agency.** The framework says nothing about intentionality, strategic choice, or the mechanisms by which actors create and sustain relationships. The bargaining models of Section 3.5 are included as examples of multi-sorted theories, but the equilibrium analysis sits inside the functor—inside the model—not in the category itself. The category records that a contest morphism exists; it does not model the decision-making process that determines its outcome.

**Theory-dependence.** The Yoneda lemma proves that identity is relational *within a given theory*. A state’s identity in the trade category may differ from its identity in the alliance category, because the two categories carry different relational structures. The framework does not provide a theory-independent notion of identity; it provides a precise notion of identity relative to a choice of relational structure. Whether this is a limitation or a feature depends on one’s philosophical commitments. For a constructivist,

the theory-dependence of identity is the point. For a positivist seeking a view from nowhere, it may be unsatisfying.

### 7.3 Extensions

Several directions suggest themselves for future work.

**Enriched categories.** The categories in this paper are “ordinary”—morphisms either exist or they don’t, with no gradations. But many IR relationships carry quantitative weight: trade volumes, alliance reliability scores, diplomatic intensity. *Enriched categories*, in which hom-sets are replaced by objects in a monoidal category (e.g., non-negative real numbers, probability distributions), can accommodate weighted relations. An enriched Yoneda lemma exists and would yield a weighted version of the relational-identity claim: identity is not just the pattern of relations but their intensities.

**Higher categories.** The paper has noted (Section 2.1) that category theory can accommodate relationships between relationships. *2-categories* and higher categorical structures formalize this: morphisms between morphisms (2-morphisms) can encode changes in relationships, renegotiations of treaties, or revisions of trade agreements. This is one route to addressing the dynamics limitation: the evolution of the relational structure can be modeled as higher-dimensional morphisms within a richer categorical framework.

**Sheaves.** The Yoneda embedding maps a category into its category of presheaves  $[\mathcal{C}^{\text{op}}, \mathbf{Set}]$ . When the category carries a notion of “covering”—a specification of which collections of relationships jointly determine a global property—the relevant models are not arbitrary presheaves but *sheaves*: presheaves satisfying a local-to-global consistency condition. Sheaf theory could formalize the idea that a state’s identity is determined not by its bilateral relationships individually but by coherent combinations of relationships that “cover” the state’s position in the system.

**Dependent type theory.** The functorial semantics used here—theories as categories, models as functors to  $\mathbf{Set}$ —is the simplest instance of a much richer correspondence between logic and category theory. Dependent type theory, which allows types to depend on terms, corresponds to categories

with richer internal structure (locally cartesian closed categories). An IR theory formalized in dependent type theory could express claims like “the nature of state  $A$ ’s alliance with state  $B$  depends on the specific trade relationship between them”—dependencies between relations that ordinary categories cannot capture. This is the most speculative extension, but it points toward a framework in which the interplay between different relational types is itself subject to formal analysis.

## 8 Conclusion

This paper asked a simple question: what is a state? The answer it offered is that a state, within any relational theory, is nothing more and nothing less than its complete relational profile—the totality of relationships that every other actor in the system bears to it.

The argument proceeded in layers. A theory of international relations is a category: a collection of actors and typed relations, with composition and identity (Section 2). A model is a functor to **Set**: a structure-preserving interpretation that assigns data to actors and functions to relations (Section 3). A translation between models is a natural transformation: a systematic conversion of one interpretation into another that respects the relational structure at every point (Section 4). The Yoneda lemma (Section 5) then proved that every datum a model assigns to a state is a global structural commitment—a complete reading of the state’s relational position, compressed into a single element—and the Yoneda embedding proved that two states are structurally identical if and only if their relational profiles agree.

These are not philosophical claims. They are theorems, valid in any category. The constructivist insight that identity is constituted by interaction, the structuralist insight that position matters, and the metatheoretic demand for precise conditions of commensurability between paradigms all find formal expression in the same framework.

Along the way, the framework delivered concrete returns. The natural-ity condition on translations between the Fearon and Beviá–Corchón bargaining models identified a sharp constraint—equilibrium war costs must be equal—that separates the region of Fearon’s parameter space that can be microfounded by a contest game from the region that cannot (Section 4.5). The analysis of existing datasets as functors revealed the implicit theoretical commitments encoded in data infrastructure: attribute tables presuppose a

discrete ontology; bilateral datasets discard composition; and each discarding has a precise cost, measured by the Yoneda lemma’s loss of bite (Section 5.5). The comparison of theories via functors gave “does trade cause peace?” a determinate form: a functor between the trade and conflict categories, subject to compositional constraints (Section 6.3).

The framework has limitations. It is static, not dynamic. It says nothing about agency or intentionality. It makes identity theory-dependent—a feature for constructivists, a cost for those seeking theory-independent foundations. Extensions to enriched categories, higher categories, and dependent type theory may address some of these limitations; others may require fundamentally different tools.

But the core result stands without qualification. In any theory that specifies actors and relations and satisfies the minimal axioms of composition, identity, and associativity, the Yoneda lemma guarantees that identity is relational. No hidden essence, no intrinsic property, no attribute not already encoded in the relational profile can distinguish two actors whose profiles agree. The state is its relations. That is not a stance. It is a theorem.

## References

- Beviá, C. and Corchón, L. C. (2010). Peace agreements without commitment. *Games and Economic Behavior*, 68(2):469–487.
- Fearon, J. D. (1995). Rationalist explanations for war. *International Organization*, 49(3):379–414.
- Fong, B. and Spivak, D. I. (2019). *An Invitation to Applied Category Theory: Seven Sketches in Compositionality*. Cambridge University Press, Cambridge.
- Hafner-Burton, E. M., Kahler, M., and Montgomery, A. H. (2009). Network analysis for international relations. *International Organization*, 63(3):559–592.
- Hoff, P. D. and Ward, M. D. (2004). Modeling dependencies in international relations networks. *Political Analysis*, 12(2):160–175.
- Jackson, P. T. (2011). *The Conduct of Inquiry in International Relations:*

- Philosophy of Science and Its Implications for the Study of World Politics.* Routledge, London.
- Kurki, M. (2008). *Causation in International Relations: Reclaiming Causal Analysis.* Cambridge University Press, Cambridge.
- Lawvere, F. W. (1963). Functorial semantics of algebraic theories. *Proceedings of the National Academy of Sciences*, 50(5):869–872.
- Mac Lane, S. (1998). *Categories for the Working Mathematician.* Springer, New York, 2nd edition.
- Maoz, Z. (2012). How network analysis can inform the study of international relations. *Conflict Management and Peace Science*, 29(3):247–256.
- Monteiro, N. P. and Ruby, K. G. (2009). IR and the false promise of philosophical foundations. *International Theory*, 1(1):15–48.
- Spivak, D. I. (2014). *Category Theory for the Sciences.* MIT Press, Cambridge, MA.
- Waltz, K. N. (1979). *Theory of International Politics.* Addison-Wesley, Reading, MA.
- Wendt, A. (1992). Anarchy is what states make of it: The social construction of power politics. *International Organization*, 46(2):391–425.
- Wendt, A. (1998). On constitution and causation in international relations. *Review of International Studies*, 24(special issue):101–117.
- Wendt, A. (1999). *Social Theory of International Politics.* Cambridge University Press, Cambridge.
- Zhukov, Y. M. and Stewart, B. M. (2013). Choosing your neighbors: Networks of diffusion in international relations. *International Studies Quarterly*, 57(2):271–287.