# Bi-interpretability from a Categorical Viewpoint

Hisashi Aratake*

*Research Institute for Mathematical Sciences, Kyoto University, Kyoto, Japan*

## 1   Introduction

First-order categorical logic (FOCL for short) originated as a categorical foundation for model theory (in Makkai & Reyes [8]). As we can see in [7], some classical model-theoretic phenomena can be efficiently described in terms of FOCL. However, most concepts in modern model theory remain to be under categorical consideration. Our principal aim is to explore categorical aspects of model theory.

This short note is intended to give an exposition for model theorists of the author's recent work on categorical characterization of bi-interpretability [1]. From the viewpoint of FOCL, classical first-order theories give rise to **Boolean pretoposes**, i.e. categories equipped with logical operations and quotients of equivalence relations. They are called **classifying pretoposes** of theories (see the definition in Theorem2.12). As Harnik [3] pointed out, a construction of classifying pretoposes can be given via Shelah's eq-construction. In fact, two theories have equivalent classifying pretoposes precisely when they are bi-interpretable in the usual model-theoretic sense. In this note, we will sketch a proof of this theorem. To keep this article accessible to model theorists, we only use elementary category theory and omit most of the details of definitions and proofs. The reader is invited to consult [1] for details.

Our theorem suggests a potential of FOCL for application to model theory. For example, we can define "model-theoretic properties" of pretoposes (e.g. completeness and stability) for those invariant under bi-interpretability. An expected application will be proposed at the end of this note.

**Terminologies and notational conventions.**   Throughout this paper, we consider many-sorted classical first-order theories. For notations and terms, we basically follow Johnstone [5, 6] (in particular, Chapter D1 in vol. 2) with several exceptions. In the following list, $\mathcal{L}$ is a fixed (many-sorted first-order) language.

- $\mathcal{L}$-Sort (resp. $\mathcal{L}$-Rel, $\mathcal{L}$-Func) denotes the set of $\mathcal{L}$-sorts (resp. relation symbols, function symbols).

---

*E-mail address*: aratake@kurims.kyoto-u.ac.jp

- If we mention finite strings of sorts or variables, the empty case is in consideration.

- A string of sorts is referred to as a **type**, and a string of variables as a **context**. The former is denoted by, say, $\bar{A}$ while the latter is denoted in bold face, say $\boldsymbol{x}$. If $\boldsymbol{x} \equiv x_1, \ldots, x_n$ where $x_i$ is a variable of sort $A_i$, we say that $\boldsymbol{x}$ is of type $\bar{A}$, and write $\boldsymbol{x} : \bar{A}$.

- We occasionally regard a 0-ary function symbol as a constant symbol.

- Let $\varphi(\boldsymbol{x})$ be an $\mathcal{L}$-formula-in-context. We often write $T \models \varphi(\boldsymbol{x})$ to mean $T \models \forall \boldsymbol{x} \varphi(\boldsymbol{x})$.

- We write $\varphi \subseteq \psi$ for $\mathcal{L}$-formulas-in-context $\varphi, \psi$, if they are in the same context (say $\boldsymbol{x} : \bar{A}$) and $T \models (\varphi(\boldsymbol{x}) \rightarrow \psi(\boldsymbol{x}))$. Although $T$ does not appear in the notation, we will always make it clear what the ambient theory is.

- $T\text{-}\mathbf{Mod}(\mathbf{Set})_e$ denotes the category of $T$-models and elementary embeddings.

## 2  Preliminaries

### 2.1  Interpretation

First, we recall the notion of interpretation in model theory while we will not pursue model-theoretic aspects of interpretations (see Hodges [4, Chap. 5 §3]). Throughout this subsection, we suppose that $\mathcal{L}, \mathcal{L}'$ are languages, and that $T$ (resp. $T'$) is an $\mathcal{L}$-theory (resp. an $\mathcal{L}'$-theory).

**Definition 2.1.** Let $\partial(\boldsymbol{x})$ be an $\mathcal{L}$-formula-in-context and $\Delta(\boldsymbol{x}, \boldsymbol{x}')$ an $\mathcal{L}$-formula-in-context with $\Delta \subseteq \partial(\boldsymbol{x}) \wedge \partial(\boldsymbol{x}')$ (where $\boldsymbol{x}, \boldsymbol{x}'$ are assumed to be disjoint). We say that $\Delta$ is a **partial $T$-equivalence relation** (on $\partial$) if the following axioms are valid in $T$:

$$\forall \boldsymbol{x} \left[\partial(\boldsymbol{x}) \rightarrow \Delta(\boldsymbol{x}, \boldsymbol{x})\right],$$
$$\forall \boldsymbol{x} \boldsymbol{x}' \left[\Delta(\boldsymbol{x}, \boldsymbol{x}') \rightarrow \Delta(\boldsymbol{x}', \boldsymbol{x})\right],$$
$$\forall \boldsymbol{x} \boldsymbol{x}' \boldsymbol{x}'' \left[\Delta(\boldsymbol{x}, \boldsymbol{x}') \wedge \Delta(\boldsymbol{x}', \boldsymbol{x}'') \rightarrow \Delta(\boldsymbol{x}, \boldsymbol{x}'')\right].$$

We say that an $\mathcal{L}$-formula $\varphi \subseteq \partial$ is **closed under** $\Delta$ when $T \models \Delta(\boldsymbol{x}, \boldsymbol{x}') \rightarrow [\varphi(\boldsymbol{x}) \leftrightarrow \varphi(\boldsymbol{x}')]$. $\qquad\square$

To define interpretations, we need a predefinition:

**Definition 2.2** (Pre-interpretation). A **pre-interpretation** $I$ of $T$ in $T'$ consists of the following data:

- For each sort $A$ in $\mathcal{L}$, we have a pair $(\partial_A^I(\boldsymbol{u}), \Delta_A^I(\boldsymbol{u}, \boldsymbol{u}'))$ where $\partial_A^I$ is an $\mathcal{L}'$-formula and $\Delta_A^I$ is a partial $T'$-equivalence relation on $\partial_A^I$.

  For each finite string of sorts $\bar{A} = A_1 \cdots A_n$, we put

  $$\partial_{\bar{A}}^I(\boldsymbol{u}) \equiv \bigwedge_{i=1}^{n} \partial_{A_i}^I(\boldsymbol{u}_i), \qquad \Delta_{\bar{A}}^I(\boldsymbol{u}, \boldsymbol{u}') \equiv \bigwedge_{i=1}^{n} \Delta_{A_i}^I(\boldsymbol{u}_i, \boldsymbol{u}_i'),$$

where all the contexts $\boldsymbol{u}_i$ and $\boldsymbol{u}'_i$ are supposed to be disjoint.

- For each relation symbol $R \rightarrowtail \bar{A}$ in $\mathcal{L}$, we have an $\mathcal{L}'$-formula $R^I(\boldsymbol{u})$ $\subseteq \partial^I_{\bar{A}}(\boldsymbol{u})$ which is closed under $\Delta^I_{\bar{A}}$.

- For each function symbol $f \colon \bar{A} \to B$ in $\mathcal{L}$, we have an $\mathcal{L}'$-formula $\Gamma^I_f(\boldsymbol{u}, \boldsymbol{v})$ $\subseteq \partial^I_{\bar{A}B}(\boldsymbol{u}, \boldsymbol{v})$ such that $\Gamma^I_f(\boldsymbol{u}, \boldsymbol{v})$ is closed under $\Delta^I_{\bar{A}B}$ and the following formulas are valid in $T'$:

$$\Delta^I_{\bar{A}}(\boldsymbol{u}, \boldsymbol{u}') \to \exists \boldsymbol{v} \left[ \Gamma^I_f(\boldsymbol{u}, \boldsymbol{v}) \wedge \Gamma^I_f(\boldsymbol{u}', \boldsymbol{v}) \right],$$
$$\exists \boldsymbol{u} \left[ \Gamma^I_f(\boldsymbol{u}, \boldsymbol{v}) \wedge \Gamma^I_f(\boldsymbol{u}, \boldsymbol{v}') \right] \to \Delta^I_B(\boldsymbol{v}, \boldsymbol{v}').$$

□

By induction on the construction of $\mathcal{L}$-formulas, we can associate each $\mathcal{L}$-formula $\varphi$ of type $\bar{A}$ with an $\mathcal{L}'$-formula $\varphi^I \subseteq \partial^I_{\bar{A}}$ which is closed under $\Delta$.

**Definition 2.3** (Interpretation). A pre-interpretation $I$ of $T$ in $T'$ is said to be an **interpretation** if $\varphi^I$ is valid in $T'$ for any $\mathcal{L}$-sentence $\varphi \in T$. When $I$ is an interpretation, we denote it by $I \colon T \to T'$. □

Let $I$ be a pre-interpretation of $T$ in $T'$ and let $\mathcal{N}$ a $T'$-model. Then we can define a canonical $\mathcal{L}$-structure $\mathcal{M}$ by using appropriate quotients $\partial^I(\mathcal{N})/\Delta^I(\mathcal{N})$. For each relation symbol $R \rightarrowtail \bar{A}$, $R^{\mathcal{M}}$ is defined to be the subset $R^I(\mathcal{N})/\Delta^I_{\bar{A}}(\mathcal{N})$ of $\partial^I_{\bar{A}}(\mathcal{N})/\Delta^I_{\bar{A}}(\mathcal{N}) = \bar{A}^{\mathcal{M}}$ [1]. Function symbols are interpreted similarly. We will denote $\mathcal{M}$ by $\mathcal{N}|_I$. The next theorem is basic:

**Theorem 2.4** (Reduction theorem, see Hodges [4, Theorem 5.3.2]). In the above notations, for any $\mathcal{L}$-formula $\varphi(\boldsymbol{x})$ and for any tuple $(\boldsymbol{a}_1, \ldots, \boldsymbol{a}_n) \in \partial^I_{\bar{A}}(\mathcal{N})$,

$$\mathcal{N} \models \varphi^I(\boldsymbol{a}_1, \ldots, \boldsymbol{a}_n) \iff \mathcal{N}|_I \models \varphi([\boldsymbol{a}_1], \ldots, [\boldsymbol{a}_n])$$

where each $[\boldsymbol{a}_i]$ is the equivalence class of $\boldsymbol{a}_i$ in $\partial^I_{A_i}(\mathcal{N})/\Delta^I_{A_i}(\mathcal{N})$. □

**Corollary 2.5.** With the same assumptions as in the Theorem, the following are equivalent:

(i) $I$ is an interpretation.

(ii) For any $T'$-model $\mathcal{N}$, $\mathcal{N}|_I$ is a $T$-model.

*Proof.* (i) $\iff T' \models \varphi^I$ holds for any $\varphi \in T$,

$\iff$ Every $T'$-model $\mathcal{N}$ satisfies $\varphi^I$ for any $\varphi \in T$,

$\iff \mathcal{N}|_I \models \varphi$ holds for any $T'$-model $\mathcal{N}$ and for any $\varphi \in T$,

$\iff$ (ii).

■

In fact, the above correspondence $\mathcal{N} \mapsto \mathcal{N}|_I$ constitutes a functor

$$(-)|_I \colon T'\text{-}\mathbf{Mod}(\mathbf{Set})_e \to T\text{-}\mathbf{Mod}(\mathbf{Set})_e. \tag{2.1}$$

We will see that this functor can be obtained by composition with a coherent functor (see (2.2)).

We now recall Shelah's eq-construction:

---

[1] Recall that $R^I(\mathcal{N})$ is closed under $\Delta^I_{\bar{A}}(\mathcal{N})$, and hence this definition makes sense.

**Definition 2.6.** We construct a language $\mathcal{L}^{\mathrm{eq}}$ and an $\mathcal{L}^{\mathrm{eq}}$-theory $T^{\mathrm{eq}}$ as follows: For each (total) $T$-equivalence relation ($\Delta$ on $\bar{A}$, say), we add to $\mathcal{L}$ a new sort $Q_\Delta$ and a new function symbol $\varepsilon_\Delta \colon \bar{A} \to Q_\Delta$, obtaining $\mathcal{L}^{\mathrm{eq}}$-Sort and $\mathcal{L}^{\mathrm{eq}}$-Func. The relation symbols of $\mathcal{L}^{\mathrm{eq}}$ are the same as $\mathcal{L}$. We then have defined $\mathcal{L}^{\mathrm{eq}}$. The $\mathcal{L}^{\mathrm{eq}}$-theory $T^{\mathrm{eq}}$ consists of the axioms of $T$ together with those of the forms

$$\forall y \exists x (\varepsilon_\Delta(x) = y) \quad \text{and} \quad \forall x \forall x'[\Delta(x, x') \leftrightarrow \varepsilon_\Delta(x) = \varepsilon_\Delta(x')].$$

□

For each partial $T$-equivalence relation $\Delta$, we associate a total $T$-equivalence relation

$$\tilde{\Delta}(x, x') \equiv (\partial(x) \wedge \partial(x') \to \Delta(x, x')) \wedge (\neg(\partial(x) \wedge \partial(x')) \to x = x').$$

Returning to our settings, for each $\mathcal{L}$-formula $\varphi(x)$, we have an $\mathcal{L}'^{\mathrm{eq}}$-formula

$$\left(\varphi^I / \Delta_{\bar{A}}^I\right)(v) \equiv \exists u \left[\varepsilon_{\tilde{\Delta}_{\bar{A}}^I}(u) = v \wedge \varphi^I(u)\right]$$

where $v$ is a variable of sort $Q_{\tilde{\Delta}_{\bar{A}}^I}$. As a result, the interpretation $I$ gives rise to a map $\varphi(x) \mapsto \left(\varphi^I / \Delta_{\bar{A}}^I\right)(v)$ from the set of all $\mathcal{L}$-formulas to that of all $\mathcal{L}'^{\mathrm{eq}}$-formulas.

## 2.2 Syntactic Category and Classifying Pretopos

**Definition 2.7** (Syntactic category)**.** For any $\mathcal{L}$-theory $T$, we construct a category $\mathcal{C}_T$ (called the **syntactic category** of $T$) as follows:

- The objects are $\alpha$-equivalence classes of $\mathcal{L}$-formulas-in-context. The $\alpha$-equivalence class of $\varphi(x)$ is denoted by $\{x.\,\varphi\}$.

- Assuming that the contexts $x, y$ are disjoint, the morphisms from $\{x.\,\varphi\}$ to $\{y.\,\psi\}$ are $T$-provably-equivalence classes of $T$-**provably functional formulas**. A $T$-provably functional formula is an $\mathcal{L}$-formula-in-context $\chi(x, y)$ such that the formulas

$$\chi \to \varphi \wedge \psi, \quad \varphi \to \exists y \chi, \quad \text{and} \quad \chi \wedge \chi[z/y] \to y = z$$

  are provable in $T$. $[\chi]$ denotes the $T$-provably-equivalence class of $\chi$.

- If $[\chi] \colon \{x.\,\varphi\} \to \{y.\,\psi\}$ and $[\theta] \colon \{y.\,\psi\} \to \{z.\,\xi\}$ are morphisms in $\mathcal{C}_T$, the composite $[\theta] \circ [\chi] \colon \{x.\,\varphi\} \to \{z.\,\xi\}$ is defined to be $[\exists y(\chi \wedge \theta)]$. □

We will follow the notational conventions below:

- We frequently write $\varphi(x)$ (or more simply $\varphi$) for an object $\{x.\,\varphi\}$.

- We often identify the object $\{x.\,x = x\}$ with the type $\bar{A}$ of $x$. We also identify the object $\{x.\,\top\}$ with the type $\bar{A}$ of $x$.

When we mention a subobject in a syntactic category, thanks to [6, Lemma D1.4.4(iv)], we may assume that it is represented by a formula in the same context.

**Proposition 2.8.** For any theory $T$, the syntactic category $\mathcal{C}_T$ is a **Boolean coherent category**, *i.e.*, roughly speaking,

- Any subobject poset $\mathrm{Sub}(\{x.\,\varphi\})$ is a Boolean algebra.

- Each morphism has an image factorization.

- These categorical structures are stable under pullbacks. $\qquad\square$

A **coherent functor** is a functor between (Boolean) coherent categories which preserves coherent structures.

**Theorem 2.9** ([6, Theorem D1.4.7]). Each $T$-model $\mathcal{M}$ gives a coherent functor $F_{\mathcal{M}}\colon \mathcal{C}_T \to \mathbf{Set}$ which sends $\{x.\,\varphi\}$ to the definable set $\varphi(\mathcal{M})$.

Moreover, this correspondence yields (a half part of) an equivalence of categories

$$T\text{-}\mathbf{Mod}(\mathbf{Set})_e \simeq \mathfrak{Coh}(\mathcal{C}_T, \mathbf{Set})$$

where $\mathfrak{Coh}(\mathcal{C}_T, \mathbf{Set})$ is the category of coherent functors from $\mathcal{C}_T$ to $\mathbf{Set}$ and natural transformations. $\qquad\square$

As we saw at the end of §2.1, an interpretation $I$ gives rise to a map $\varphi(x) \mapsto \left(\varphi^I/\Delta_{\bar{A}}^I\right)(v)$ from the set of all $\mathcal{L}$-formulas to that of all $\mathcal{L}'^{\mathrm{eq}}$-formulas. Now we can associate $I$ with a coherent functor:

**Proposition 2.10.** An interpretation $I\colon T \to T'$ induces a coherent functor $F_I\colon \mathcal{C}_T \to \mathcal{C}_{T'^{\mathrm{eq}}}$ sending a morphism $[\chi]\colon \{x.\,\varphi\} \to \{y.\,\psi\}$ in $\mathcal{C}_T$ to the following morphism in $\mathcal{C}_{T'^{\mathrm{eq}}}$:

$$\left[\exists u \exists v \left(\chi^I(u,v) \wedge \varepsilon_{\tilde{\Delta}_{\bar{A}}^I}(u) = u' \wedge \varepsilon_{\tilde{\Delta}_{\bar{B}}^I}(v) = v'\right)\right] \colon \{u'.\,\varphi^I/\Delta_{\bar{A}}^I\} \to \{v'.\,\psi^I/\Delta_{\bar{B}}^I\}.$$

$$\square$$

We can make up $F_I\colon \mathcal{C}_T \to \mathcal{C}_{T'^{\mathrm{eq}}}$ into another coherent functor $\mathcal{P}_I\colon \mathcal{C}_{T^{\mathrm{eq}}} \to \mathcal{C}_{T'^{\mathrm{eq}}}$. The category $\mathcal{C}_{T^{\mathrm{eq}}}$ will play a special role in our theory.

**Definition 2.11** (Proper theory). A classical theory $T$ is said to be **proper** if there exists a sort $D$ such that $T \models \exists x x'(x \neq x')$ with $x, x' : D$. $\qquad\square$

In what follows, all theories are assumed to be proper. The following observation was made by Harnik [3].

**Theorem 2.12.** Let $T$ be a theory. Then $\mathcal{C}_{T^{\mathrm{eq}}}$ is a **Boolean pretopos**, i.e. a Boolean coherent category having the following properties:

- Any "equivalence relation" has a quotient.

- Finite coproducts exist and are disjoint.

Moreover, the canonical functor $\gamma\colon \mathcal{C}_T \to \mathcal{C}_{T^{\mathrm{eq}}}$ gives a **pretopos completion** of $\mathcal{C}_T$: for any (Boolean) pretopos $\mathcal{P}$ and for any coherent functor $F\colon \mathcal{C}_T \to \mathcal{P}$,

there exists a unique coherent functor $G\colon \mathcal{C}_{T^{\mathrm{eq}}} \to \mathcal{P}$ (up to natural isomorphism) such that $F \simeq G\gamma$:

$$
\begin{array}{ccc}
\mathcal{C}_T & \xrightarrow{\;\gamma\;} & \mathcal{C}_{T^{\mathrm{eq}}} \\
& F\searrow & \downarrow G \\
& & \mathcal{P}
\end{array}
$$

In fact, the following equivalences hold:

$$T\text{-}\mathbf{Mod}(\mathbf{Set})_e \simeq \mathfrak{Coh}(\mathcal{C}_T, \mathbf{Set}) \simeq \mathfrak{Coh}(\mathcal{C}_{T^{\mathrm{eq}}}, \mathbf{Set}) \simeq T^{\mathrm{eq}}\text{-}\mathbf{Mod}(\mathbf{Set})_e.$$

The **classifying pretopos** $\mathcal{P}_T$ of $T$ is defined to be $\mathcal{C}_{T^{\mathrm{eq}}}$. $\square$

**Corollary 2.13.** Combining Proposition 2.10 and Theorem 2.12, we can obtain from an interpretation $I\colon T \to T'$ another coherent functor $\mathcal{P}_I\colon \mathcal{P}_T \to \mathcal{P}_{T'}$. $\square$

Observe that the restriction functor $(-)|_I\colon T'\text{-}\mathbf{Mod}(\mathbf{Set})_e \to T\text{-}\mathbf{Mod}(\mathbf{Set})_e$ which appeared in (2.1) can be described as the composite of the functors below:

$$
\begin{array}{ccccccc}
\mathfrak{Coh}(\mathcal{C}_{T'}, \mathbf{Set}) & \xrightarrow{\sim} & \mathfrak{Coh}(\mathcal{P}_{T'}, \mathbf{Set}) & \xrightarrow{(-)\circ\mathcal{P}_I} & \mathfrak{Coh}(\mathcal{P}_T, \mathbf{Set}) & \xrightarrow{\sim} & \mathfrak{Coh}(\mathcal{C}_T, \mathbf{Set}) \\
\downarrow\wr & & & & & & \downarrow\wr \\
T'\text{-}\mathbf{Mod}(\mathbf{Set})_e & & & \xrightarrow{\quad (-)|_I \quad} & & & T\text{-}\mathbf{Mod}(\mathbf{Set})_e
\end{array}
$$

$$(2.2)$$

# 3   Bi-interpretability

Let $I$ and $J$ be interpretations of $T$ in $T'$.

**Definition 3.1** (Homotopy between interpretations)**.** A **homotopy** $h\colon I \Rightarrow J$ is a family $\{h_A\}_A$ of $\mathcal{L}'$-formulas with $h_A \subseteq \partial_A^I \wedge \partial_A^J$ such that

- Each $h_A$ is an $\mathcal{L}'$-formula closed under $\Delta_A^I \wedge \Delta_A^J$.

- Each $h_A$ induces an isomorphism $\tilde{h}_A\colon \partial_A^I/\Delta_A^I \xrightarrow{\sim} \partial_A^J/\Delta_A^J$ in $\mathcal{P}_{T'}$ similarly to the definition of $F_I([\chi])$ in Proposition 2.10.

- For any atomic $\mathcal{L}$-formula $\varphi(\boldsymbol{x})$, these isomorphisms induce an isomorphism $\varphi^I/\Delta_{\bar{A}}^I \xrightarrow{\sim} \varphi^J/\Delta_{\bar{A}}^J$.

If we have two homotopies $h, k\colon I \Rightarrow J$ and moreover if their components $h_A$ and $k_A$ are $T'$-equivalent for any $A$, then they will be identified.

We say that $I$ and $J$ are **homotopic** when there exists a homotopy from $I$ to $J$. This is an equivalence relation on interpretations of $T$ in $T'$. $\square$

**Proposition 3.2.** A homotopy $h\colon I \Rightarrow J$ induces a natural isomorphism $\mathcal{P}_h\colon \mathcal{P}_I \overset{\sim}{\Rightarrow} \mathcal{P}_J$.

$$
\begin{array}{ccccccc}
& I & & & & \mathcal{P}_I & \\
T & \overset{\Downarrow h}{\underset{J}{\rightrightarrows}} & T' & \rightsquigarrow & \mathcal{P}_T & \overset{\Downarrow \mathcal{P}_h}{\underset{\mathcal{P}_J}{\rightrightarrows}} & \mathcal{P}_{T'}
\end{array}
$$

□

**Definition 3.3** (Composition of interpretations). Let $I \colon T \to T'$ and $J \colon T' \to T''$ be interpretations. Define an interpretation $JI \colon T \to T''$ as follows:

- For each sort $A$ in $\mathcal{L}$, we put $\partial_A^{JI} \equiv (\partial_A^I)^J$ and $\Delta_A^{JI} \equiv (\Delta_A^I)^J$.

- For each relation symbol $R$ in $\mathcal{L}$, we put $R^{JI} \equiv (R^I)^J$.

- For each function symbol $f$ in $\mathcal{L}$, we put $\Gamma_f^{JI} \equiv (\Gamma_f^I)^J$.

These data indeed satisfy the conditions to be an interpretation.   □

**Definition 3.4** (Bi-interpretability of theories). We say that two theories $T$ and $T'$ are **bi-interpretable** when there exist two interpretations $I \colon T \to T'$ and $J \colon T' \to T$ such that

- $JI$ is homotopic to $I_T \colon T \to T$ (the identity interpretation),

- $IJ$ is homotopic to $I_{T'} \colon T' \to T'$.   □

We now describe our main result on bi-interpretability.

**Theorem 3.5** (Categorical characterization of bi-interpretability). *$T$ and $T'$ are bi-interpretable precisely when their classifying pretoposes are equivalent, i.e. $\mathcal{P}_T \simeq \mathcal{P}_{T'}$.*   □

Two proofs of this theorem are provided in [1]. One uses *two-dimensional* category theory, and we do not present its method here. The other involves the notions of *Morita extension* and *Morita span*, which are recently introduced by Barrett & Halvorson [2]. They introduced these to give a plausible definition of theoretical equivalence which generalizes definitional equivalence. The equivalence is defined by existence of a Morita span (see §3.1).

Our proof shows that these three conditions, i.e. bi-interpretability, equivalence of classifying pretoposes and existence of a Morita span, are all equivalent. In §3.2, we will briefly describe the latter proof.

## 3.1   Morita Extension

Let $\mathcal{L}$ and $\mathcal{L}^+$ be languages with $\mathcal{L} \subseteq \mathcal{L}^+$. We put for convenience

$$(\mathcal{L}^+ \setminus \mathcal{L})\text{-Sort} = \mathcal{L}^+\text{-Sort} \setminus \mathcal{L}\text{-Sort}, \qquad (\mathcal{L}^+ \setminus \mathcal{L})\text{-Rel} = \mathcal{L}^+\text{-Rel} \setminus \mathcal{L}\text{-Rel},$$
$$(\mathcal{L}^+ \setminus \mathcal{L})\text{-Func} = \mathcal{L}^+\text{-Func} \setminus \mathcal{L}\text{-Func}.$$

**Definition 3.6** (Explicit definitions).

(1) Suppose that $R : \bar{A} \in (\mathcal{L}^+ \setminus \mathcal{L})\text{-Rel}$. An **explicit definition of $R$ in terms of $\mathcal{L}$** is an $\mathcal{L}^+$-sentence of the form

$$\forall \boldsymbol{x}(R(\boldsymbol{x}) \leftrightarrow \varphi(\boldsymbol{x})),$$

where $\varphi(\boldsymbol{x})$ is an $\mathcal{L}$-formula in the same canonical context as $R$.

(2) Suppose that $f\colon \bar{A} \to B \in (\mathcal{L}^+ \setminus \mathcal{L})$-Func. An **explicit definition of $f$ in terms of $\mathcal{L}$** is an $\mathcal{L}^+$-sentence of the form

$$\forall \boldsymbol{x} \forall y (f(\boldsymbol{x}) = y \leftrightarrow \varphi(\boldsymbol{x}, y)),$$

where $\varphi(\boldsymbol{x}, y)$ is an $\mathcal{L}$-formula in the same canonical context as $f$. For an arbitrary explicit definition of a function symbol, we refer to the $\mathcal{L}$-sentence

$$\forall \boldsymbol{x} \exists! y \varphi(\boldsymbol{x}, y)$$

as the **admissibility condition** for the explicit definition. If $f$ is actually a constant symbol, i.e. of arity zero, the above formulas have the following forms:

$$\forall y (y = c \leftrightarrow \varphi(y)),$$
$$\exists! y \varphi(y).$$

$\square$

**Definition 3.7** (Sort definitions). Suppose that $S \in (\mathcal{L}^+ \setminus \mathcal{L})$-Sort, $A \in \mathcal{L}$-Sort and $\bar{A}$ is a type in $\mathcal{L}$. In the following, $s$ (resp. $x, \boldsymbol{x}$) is a variable (or a context) of the sort $S$ (resp. $A, \bar{A}$).

(1) A **sort definition of $S$ as a product sort** in terms of $\mathcal{L}$ is an $\mathcal{L}^+$-sentence of the form

$$\forall \boldsymbol{x} \exists! s \bigwedge_i \pi_i(s) = x_i,$$

where all $\pi_i \colon S \to A_i$ are in $(\mathcal{L}^+ \setminus \mathcal{L})$-Func.

(2) A **sort definition of $S$ as a coproduct sort** in terms of $\mathcal{L}$ is an $\mathcal{L}^+$-sentence of the form

$$\forall s \left[ \bigvee_i \exists! x_i (\rho_i(x_i) = s) \right] \wedge \bigwedge_{i \neq j} \forall x_i \forall x_j (\rho_i(x_i) \neq \rho_j(x_j)),$$

where all $\rho_i \colon A_i \to S$ are in $(\mathcal{L}^+ \setminus \mathcal{L})$-Func.

(3) A **sort definition of $S$ as a subsort** in terms of $\mathcal{L}$ is an $\mathcal{L}^+$-sentence of the form

$$\forall x \left[ \varphi(x) \leftrightarrow \exists s (\iota_\varphi(s) = x) \right] \wedge \forall s \forall s' \left[ \iota_\varphi(s) = \iota_\varphi(s') \to s = s' \right],$$

where $\varphi(x)$ is an $\mathcal{L}$-formula-in-context and $\iota_\varphi \colon S \to A$ is in $(\mathcal{L}^+ \setminus \mathcal{L})$-Func.

(4) A **sort definition of $S$ as a quotient sort** in terms of $\mathcal{L}$ is an $\mathcal{L}^+$-sentence of the form

$$\forall x \forall x' \left[ \varepsilon_\Delta(x) = \varepsilon_\Delta(x') \leftrightarrow \Delta(x, x') \right] \wedge \forall s \exists x \left[ \varepsilon_\Delta(x) = s \right],$$

where $\Delta(x, x')$ is an $\mathcal{L}$-formula-in-context and $\varepsilon_\Delta \colon A \to S$ is in $(\mathcal{L}^+ \setminus \mathcal{L})$-Func. The **admissibility condition** for the definition is the $\mathcal{L}$-sentence which expresses that $\Delta$ is an equivalence relation on the sort $A$. $\square$

**Definition 3.8** (Morita extensions, see Barrett & Halvorson [2])**.** Let $T$ be an $\mathcal{L}$-theory and $T^+$ an $\mathcal{L}^+$-theory with $T \subseteq T^+$. We say that $T^+$ is a **Morita extension** of $T$ if the following conditions hold:
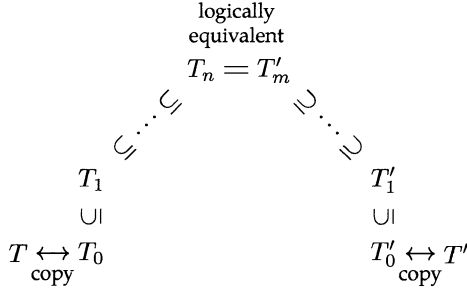
- Every $S \in (\mathcal{L}^+ \setminus \mathcal{L})$-Sort is defined as either a product sort $\prod_j A_j$, a coproduct sort $\coprod_j A_j$, a subsort $A_\varphi$ or a quotient sort $A/\Delta$ in terms of $\mathcal{L}$. This means that $(\mathcal{L}^+ \setminus \mathcal{L})$-Func contains the involved function symbols $\pi_i, \rho_i, \iota_\varphi, \varepsilon_\Delta$, that $T^+$ contains the involved sort definitions and that the involved admissibility conditions for quotient sorts are valid in $T$.

- Every $f \in (\mathcal{L}^+ \setminus \mathcal{L})$-Func other than the above $\pi_i, \rho_i, \iota_\varphi, \varepsilon_\Delta$ is explicitly defined in terms of $\mathcal{L}$. This means that $T^+$ contains the involved explicit definitions and that the involved admissibility conditions are valid in $T$.

- Every $R \in (\mathcal{L}^+ \setminus \mathcal{L})$-Rel is explicitly defined in terms of $\mathcal{L}$. This means that $T^+$ contains the involved explicit definitions.

- $T^+ \setminus T$ contains no axioms other than those mentioned above. $\qquad \square$

Morita extension is a generalized notion of definitional extension which admits sort definitions. In contrast to eq-construction, here we are allowed to define a new sort by using only sorts (not types). On the other hand, Morita extension has some similar properties with eq-construction (e.g. it is conservative).

**Morita span.** We say that an $\mathcal{L}'$-theory $T'$ is a **copy** of an $\mathcal{L}$-theory $T$ if $\mathcal{L}'$ is obtained from $\mathcal{L}$ by renaming some (possibly no) symbols in $\mathcal{L}$ to symbols not contained in $\mathcal{L}$ and $T'$ is obtained by replacing these symbols appearing in the axioms in $T$ with the corresponding symbols in $\mathcal{L}'$.

**Definition 3.9** (Morita span)**.** Let $T$ (resp. $T'$) be an $\mathcal{L}$-theory (resp. an $\mathcal{L}'$-theory). A **Morita span** from $T$ to $T'$ is a family of finitely many theories $T_0, \ldots, T_n, T'_0, \ldots, T'_m$ (possibly $n, m = 0$) which satisfy the following conditions:

- $T_i$ is an $\mathcal{L}_i$-theory for each $i = 0, \ldots, n$, and $T'_j$ is an $\mathcal{L}'_j$-theory for each $j = 0, \ldots, m$.

- $T_{i+1}$ is a Morita extension of $T_i$ for each $i = 0, \ldots, n-1$, and $T'_{j+1}$ is a Morita extension of $T'_j$ for each $j = 0, \ldots, m-1$.

- $T_0$ is a copy of $T$, and $T'_0$ is a copy of $T'$.

- $\mathcal{L}_n$ and $\mathcal{L}'_m$ are identical and $T_n$ and $T'_m$ are **logically equivalent**, i.e. they make the same sentences valid.

$$\text{logically equivalent}$$
$$T_n = T'_m$$

$$\subseteq \cdots \subseteq \qquad \supseteq \cdots \supseteq$$

$$T_1 \qquad\qquad\qquad T'_1$$
$$\cup\vert \qquad\qquad\qquad \cup\vert$$
$$T \leftrightarrow T_0 \qquad\qquad\qquad T'_0 \leftrightarrow T'$$
$$\text{copy} \qquad\qquad\qquad\qquad \text{copy}$$

$\square$

A proof for existence of a Morita span to define an equivalence relation on theories is postponed to Corollary 3.11.

## 3.2 Sketch of Proof

Notice two trivial facts:

- Any copy of $T$ is bi-interpretable with $T$.

- Two bi-interpretable theories have equivalent classifying pretoposes, for Proposition 3.2 ensures that a bi-interpretation gives an equivalence between the classifying pretoposes.

We now sketch a proof of Theorem 3.5. The proof proceeds by showing the following implications:

$$\text{Bi-interpretability} \implies \text{Equivalence of classifying pretoposes}$$
$$\implies \text{Existence of a Morita span} \implies \text{Bi-interpretability}$$

The first implication was explained above. For the third implication, the general case reduces to a Morita extension:

**Theorem 3.10.** Any two theories constituting a Morita extension are bi-interpretable. As a consequence, so are any two theories which can be connected by a Morita span.

*Proof.* Let $T^+$ be a Morita extension of $T$. We have to get an interpretation $J \colon T^+ \to T$. For $\prod_i A_i, A_\varphi, A/\Delta \in (\mathcal{L}^+ \setminus \mathcal{L})$-Sort, their interpretations under $J$ are defined as expected. For the coproduct case, using properness, we can also take an $\mathcal{L}^{\text{eq}}$-formula defining a desired coproduct in $\mathcal{P}_T$. Therefore, we have obtained a pre-interpretation $J$. This is indeed an interpretation and, together with the canonical interpretation $I \colon T \to T^+$, constitutes a bi-interpretation. ∎

The second implication needs substantial works. To obtain a desired Morita span from $T$ to $T'$, we construct sequences of Morita extensions

$$T \subseteq T_1 \subseteq \cdots \subseteq T_5 \quad \text{and} \quad T' \subseteq T'_1 \subseteq \cdots \subseteq T'_5,$$

where $T_i$ (resp $T'_j$) is an $\mathcal{L}_i$-theory (resp. an $\mathcal{L}'_j$-theory) and the following conditions hold:

- There is a bijection between $\mathcal{L}_5$ and $\mathcal{L}_5'$.

- If we identify $\mathcal{L}_5$ with $\mathcal{L}_5'$ by using the bijection, $T_5$ and $T_5'$ are logically equivalent.

The essential idea is as follows: if we have a categorical equivalence $\mathcal{P}_T \simeq \mathcal{P}_{T'}$ and skeletons $\mathcal{S}$, $\mathcal{S}'$ (i.e. complete systems of representatives of isomorphism classes of objects) of $\mathcal{P}_T$, $\mathcal{P}_{T'}$ respectively, then a bijection between $\mathcal{S}$ and $\mathcal{S}'$ is induced. We add to $T$ (resp. $T'$) some product sorts, subsorts and quotient sorts step by step such that, for each formula $\psi$ in $\mathcal{S}$ (resp. in $\mathcal{S}'$), there exists a sort in $\mathcal{L}_3$ (resp. in $\mathcal{L}_3'$) representing $\psi$. Slightly modifying $T_3$, $T_3'$, we obtain $T_4$, $T_4'$ with $\mathcal{L}_4$-Sort $\overset{\sim}{\to} \mathcal{L}_4'$-Sort. At last we add to $T_4$, $T_4'$ explicit definitions of relations and functions such that the above bijections extends to $\mathcal{L}_5 \overset{\sim}{\to} \mathcal{L}_5'$.

These constructions indeed give logically equivalent $T_5$, $T_5'$ while we suppressed lots of technical details. Once we have proved Theorem 3.5, we also fulfill a promise made at the end of §3.1.

**Corollary 3.11.** Existence of a Morita span defines an equivalence relation on theories. □

Since ($\lambda$-)stability is invariant under bi-interpretability, we have the following somewhat new definition:

**Definition 3.12** (Stability of pretoposes). A pretopos $\mathcal{P}$ is ($\lambda$-)**stable** when it is (equivalent to) the classifying pretopos of some ($\lambda$-)stable theory. □

It is known that completeness of a theory can be described categorically via the notion of *two-valued pretopos*. On the other hand, we do not know how stability of a pretopos can be described in the language of pretoposes. If we succeed in characterizing stability categorically, this will suggest more extensive uses of category theory in modern model theory. For example, we expect that 2-categorical constructions for pretoposes will give new model-theoretic constructions for good theories. So we will pursue this direction.

# References

[1]  H. Aratake. "Bicategory of Theories and Interpretations". In preparation.

[2]  T. W. Barrett and H. Halvorson. "Morita Equivalence". In: **The Review of Symbolic Logic 9**.3 (2016), pp. 556–582. arXiv: 1506.04675 [math.LO].

[3]  V. Harnik. "Model Theory vs. Categorical Logic: Two Approaches to Pretopos Completion (a.k.a. $T^{\text{eq}}$)". In: **Models, Logics, and Higher-Dimensional Categories: A Tribute to the Work of Mihály Makkai**. Ed. by B. Hart et al. CRM Proceedings & Lecture Notes 53. American Mathematical Society, 2011, pp. 79–106.

[4]  W. Hodges. **Model Theory**. Encyclopedia of Mathematics and its Applications 42. Cambridge University Press, 1993.

[5]  P. T. Johnstone. **Sketches of an Elephant: A Topos Theory Compendium**. Vol. 1. Oxford Logic Guides 43. Clarendon Press, 2002.

[6]   P. T. Johnstone. **Sketches of an Elephant: A Topos Theory Compendium**. Vol. 2. Oxford Logic Guides 44. Clarendon Press, 2002.

[7]   M. Makkai. "Duality and Definability in First Order Logic". In: **Memoirs of the American Mathematical Society 105**.503 (1993). x+106 pp.

[8]   M. Makkai and G. E. Reyes. **First Order Categorical Logic: Model-Theoretical Methods in the Theory of Topoi and Related Categories**. Lecture Notes in Mathematics 611. Springer-Verlag, 1977.