

Frege's Concept Horse Problem in the Simply-Typed λ -calculus

David Titarenco (group: Jon Ben & Jose Trujillo)

University of California, Los Angeles

Abstract

In *On Function and Concept*, G. Frege makes a distinction between concepts and objects. Specifically, objects are *saturated object-expressions* while concepts are *unsaturated function-expressions*. The goal of this paper is to illuminate what the *object-concept* distinction entails in the context of **simply-typed lambda calculus** (λ^{\rightarrow}) and whether or not this entailment is susceptible to problems like the concept horse paradox.

I. INTRODUCTION

Frege's *On Function and Concept* (1892) claims that there is a substantive difference between concepts and objects. This distinction is described on pg. 147:

Statements in general, just like equations or inequalities or expressions in Analysis, can be imagined to be split up into two parts; one complete in itself, and the other in need of supplementation, or 'unsaturated'. Thus, e.g., we split up the sentence 'Caesar conquered Gaul' into 'Caesar' and 'conquered Gaul'.

Definition 1. Complete parts of statements are 'saturated' while incomplete ones are 'unsaturated'.

Example 1. Consider the following statement:

'Fido jumped over the fence' (1)

By **Definition 1**, we can see that (1) is split up into the saturated 'Fido' object-expression and 'jumped over the fence' unsaturated function-expression.

Remark 1. These two categories, as it turns out, have simple analogues in λ^\rightarrow . Unsaturated function-expressions are terms with functional types, e.g. $N^{\sigma \rightarrow \tau}$ or $N : \sigma \rightarrow \tau$, while saturated object-expressions are terms with non-functional types, e.g. M^σ or $M : \sigma$.

Definition 2. We say that if M gets type σ and N gets type $\sigma \rightarrow \tau$ then the application NM is **legal** (as N is considered a function from terms of type σ to terms of type τ) and gets type τ .

Definition 3. The set of types of λ^\rightarrow , denoted by $\text{Type}(\lambda^\rightarrow)$ is inductively defined as follows. We write $\mathbb{T} = \text{Type}(\lambda^\rightarrow)$ where

$$\begin{aligned} \alpha, \alpha', \alpha'', \dots &\in \mathbb{T} && \text{(type variables);} \\ \sigma, \tau \in \mathbb{T} &\Rightarrow (\sigma \rightarrow \tau) \in \mathbb{T} && \text{(function space variables).} \end{aligned}$$

Definition 4. We will now define a new variant of \mathbb{T} that more closely models Frege's natural language notions of *concept* and *object*. This new variant, $\mathbb{T}_{\mathcal{F}}$, has the following restrictions:

- (a) There are three and *only* three types:
$$\begin{cases} \alpha \in \mathbb{T}_{\mathcal{F}} & (\text{objects}); \\ \mathbf{H} \in \mathbb{T}_{\mathcal{F}} & (\text{truth-values}); \\ (\alpha \rightarrow \mathbf{H}) \in \mathbb{T}_{\mathcal{F}} & (\text{concepts}). \end{cases}$$
- (b) *Objects* have an unrestricted domain.
- (c) A truth-value is a boolean type where $\mathbf{H} = \{\text{'true'}, \text{'false'}\}$.
- (d) A *concept* is a functional type that takes an *object* as an input and returns a truth value. This mapping is done by some valuation function v where v maps to ‘true’ if the function-expression applied to the argument(s) is true, e.g., in (1), if Fido *did* jump over a fence then $v(\text{'Fido'}) = \text{'true'}$; otherwise, v maps to ‘false’.
- (e) Note that $(\mathbf{H} \rightarrow \alpha) \notin \mathbb{T}_{\mathcal{F}}$.

Remark 2. Until otherwise noted, we will work only with types in $\mathbb{T}_{\mathcal{F}}$.

Definition 5. Let M be an untyped λ^{\rightarrow} -term. Given the non-functional type α , we say that $M : \alpha$ is a *saturated object-expression*, or simply an *object*.

Definition 6. Let N be an untyped λ^{\rightarrow} -term. Given the functional type $\alpha \rightarrow \mathbf{H}$, we say that $N : \alpha \rightarrow \mathbf{H}$ is an *unsaturated function-expression*, or simply a *concept*.

Definition 7. We will now make a semantic distinction between two kinds of uses of the token ‘is’ in natural language. To illustrate the distinction, consider:

- (a) ‘Mark Twain is Samuel Clemens’: here, we have the ‘*is*’ of *identity*, or $x = x$. The equality operator $=$ is used as a mathematical formalism and may not necessarily be computable in λ^{\rightarrow} . So, $v(\text{'Mark Twain'}) = \text{'true'}$ iff ‘Mark Twain’ = ‘Samuel Clemens’.
- (b) ‘Mark Twain is dead’: here, we have the ‘*is*’ of *predication* which, if combined with some property P will yield a functional type of form $\alpha \rightarrow \mathbf{H}$ for some object α . As in **Definition 4**, if α has property P , $v(\alpha) = \text{'true'}$; otherwise, v maps to ‘false’.

II. SIMPLE NATURAL LANGUAGE SENTENCES IN λ^\rightarrow

Theorem 1. Consider the same sentence from **Example 1**:

$$\text{‘Fido jumped over the fence’} \quad (2)$$

This sentence can be formulated in λ^\rightarrow like so:

$$\Gamma \vdash J : \alpha \rightarrow \mathbf{H}, \Gamma \vdash F : \alpha \Rightarrow \Gamma \vdash (JF) : \mathbf{H} \quad (3)$$

Proof. We will set the *object-expression* ‘Fido’ = F and the (one-place) *function-expression* ‘ (x) jumped over the fence’ = $J(x)$. Note that we will drop the (x) in practice for λ -calculus wff satisfiability. By **Definition 5**, F must be of type α . And to make (JF) legal, by **Definition 2**, J must be of type $\alpha \rightarrow \mathbf{H}$. Therefore, a correctly typed (JF) from $J : \alpha \rightarrow \mathbf{H}$ and $F : \alpha$ looks like

$$(JF) : \mathbf{H} \quad (4)$$

which is precisely what we wanted to show in (3). \square

Proof. We can also derive a more formally-rigorous proof. Consider:

$$\frac{\Gamma \vdash J : \alpha \rightarrow \mathbf{H} \quad \Gamma \vdash F : \alpha}{\Gamma \vdash JF : \mathbf{H}} \text{ } (\rightarrow e) \quad (5)$$

where we get the same result: JF is of type \mathbf{H} . \square

Remark 3. Natural language predicates (e.g. ‘is blue’, ‘jumped over the fence’, ‘is a concept’, ‘conquered Gaul’, ‘may have stolen my wallet’, etc.) are *concepts* in $\mathbb{T}_{\mathcal{F}}$ while subjects, be they proper nouns, improper nouns, pronouns, etc. (e.g. ‘the car’, ‘Fido’, ‘he’, ‘the tall man’, ‘the Brooklyn Bridge’, etc.) are *objects*. We will now look at more nontrivial examples of natural language propositions and attempt to model them in λ^\rightarrow under $\mathbb{T}_{\mathcal{F}}$.

Example 2. Consider the statement:

$$\text{‘The concept horse is a concept.’} \quad (6)$$

Claim 1. We will set the *object-expression* ‘The concept horse’ = T_{ch} and the (one-place) *function-expression* ‘ (x) is a concept’ = C . **Example 2** is witness to a semantic ambiguity between natural language and our type system ($\mathbb{T}_{\mathcal{F}}$). More specifically, (6) evaluates to

‘true’ in natural language, but ‘false’ under $\mathbb{T}_{\mathcal{F}}$, i.e. $CT_{ch} : \mathbf{H}$ has an ambiguous truth-value. So, if $T_{ch} : \alpha$ is the *saturated object-expression* ‘the concept horse’ and $C : \alpha \rightarrow \mathbf{H}$ is the *unsaturated function-expression* ‘is a concept’, then the mapping $v : \alpha \rightarrow \mathbf{H}$ is ambiguous.

Proof. Consider the following natural language statement:

$$\text{‘The house is a house.’} \tag{7}$$

Trivially, (7) is true. Per **Definition 7**, the ‘is’ here seems to be an ‘*is*’ of *identity*, so $v(\text{‘house’}) = \text{‘true’}$ since ‘house’ = ‘house’.

Now, let’s define a new set of types $\mathbb{T}_{\mathcal{F}}^*$ where $\mathbb{T}_{\mathcal{F}}^*$ is just like $\mathbb{T}_{\mathcal{F}}$ except that $\mathbb{T}_{\mathcal{F}}^*$ has one and only one additional type: $\beta \in \mathbb{T}_{\mathcal{F}}^*$. We will call all terms of type β *houses*. Now, we an ambiguity in (7); it is unclear whether we’re asking if the identity ‘house’ = ‘house’ holds or whether we’re asking if the *saturated object-expression* ‘the house’ is of type β . The former is true per the first half of this proof. The latter is false, as any *object-expression* is of type α in $\mathbb{T}_{\mathcal{F}}$ and, by extension, also in $\mathbb{T}_{\mathcal{F}}^*$ and $\alpha \neq \beta$. This same kind of ambiguity happens in (7). \square

Definition 8. Given the proof of *Claim 1*, we will now denote two variations of v , essentially defining two valuations for $\alpha \rightarrow \mathbf{H}$:

- (a) v_0 for meta-language predicates;
- (b) $v_{\mathcal{L}}$ for natural language predicates.

In light of this, we now have two interpretations for statements that involve *objects*, *concepts*, and truth-values. So, in some cases, $v_0(CT_{ch} : \mathbf{H})$ will evaluate to ‘false’, but $v_{\mathcal{L}}(CT_{ch} : \mathbf{H})$ evaluates to ‘true’ or vice-versa. This is known as the *concept horse paradox*.

Theorem 2. Given any object-expression, $\mathbb{O} : \alpha$ and if we let the function-expression ‘(x) is a concept’ = C , $v_0(C\mathbb{O} : \mathbf{H})$ will always evaluate to ‘false’.

Proof. From the proof of *Claim 1* and **Definition 8**. \square

Example 3. Now consider the statement:

$$\text{‘Seabiscuit is the concept horse.’} \tag{8}$$

Claim 2. We will set the *object-expression* ‘Seabiscuit’ = S and the (one-place) *function-expression* ‘(x) is the concept horse’ = C_h . **Example 3** avoids the ambiguity in (6) because its truth value happens to be false under both valuations, so $v_0(C_h S : \mathbf{H}) = v_{\mathfrak{L}}(C_h S : \mathbf{H})$.

Proof. By **Theorem 2**, $v_0(C_h S : \mathbf{H})$, will always evaluate to ‘false’. In natural language, ‘Seabiscuit’ is not ‘the concept horse’ (at best, it is an instance of ‘the concept horse’), so $v_{\mathfrak{L}}(C_h S : \mathbf{H})$ will evaluate to ‘false’ here as well. Since both valuations return ‘false’, we avoid the paradox. \square

Example 4. Consider the statement:

$$\text{‘Anna is something Bill is not – namely, a student.’} \quad (9)$$

Remark 4. **Example 4** is not ambiguous as it only uses terms found in natural language. This is trivial to prove (since the only valid valuation is $v_{\mathfrak{L}}$) However, (9) introduces another problem with $\mathbb{T}_{\mathcal{F}}$ (more specifically, with λ^{\rightarrow}): subtyping. It seems intuitive that ‘being a student’ should be a member of ‘things Bill is not’. As it turns out, λ^{\rightarrow} has no way of representing this.

III. PROPOSED SOLUTIONS

A. Semantic Culling

The most straightforward solution to the semantic problem of the concept horse paradox is a simple one. Instead of using terms like *concept* and *object*, that not only have a meaning in the meta-language, but also in natural languages, we will use terms that have no meaning in natural language. Consider a new type system, $T_{\mathcal{F}}^{**}$.

$$\text{The types of } T_{\mathcal{F}}^{**}: \begin{cases} \alpha \in T_{\mathcal{F}}^{**} & (\text{foo}); \\ \mathbf{H} \in T_{\mathcal{F}}^{**} & (\text{baz}); \\ (\alpha \rightarrow \mathbf{H}) \in T_{\mathcal{F}}^{**} & (\text{foobaz}). \end{cases}$$

Under such a schema, statements like ‘The concept horse is [a] foobaz’ are meaningless in natural language, but are semantically well-formed in the meta-language. Thus, we can infer that we need to evaluate the statement with v_0 as opposed to $v_{\mathfrak{L}}$ and avoid any ambiguity.

Remark 5. The caveat here is that we lose some expressiveness in natural languages. In **Example 3**, we were able to refer both to the semantics of natural language and the meta-language without ambiguity. This would no longer be possible.

B. Syntactic Sugar

Another proposed solution is syntactic in nature. We previously showed that ‘Fido jumped over the fence’ (**Example 1**) can be formulated as follows:

$$\frac{\Gamma \vdash J : \alpha \rightarrow \mathbf{H} \quad \Gamma \vdash F : \alpha}{\Gamma \vdash JF : \mathbf{H}} \quad (\rightarrow e) \quad (10)$$

We will now introduce a new type of λ^{\rightarrow} -term that will differentiate between what valuation v one should use to determine the truth-value of some $M : \mathbf{H}$. To do this, consider a new type system $\mathbb{T}_{\mathcal{F}}^{\omega}$.

$$\text{The types of } \mathbb{T}_{\mathcal{F}}^{\omega}: \begin{cases} \alpha \in \mathbb{T}_{\mathcal{F}}^{\omega} & (\text{objects}); \\ \mathbf{H} \in \mathbb{T}_{\mathcal{F}}^{\omega} & (\text{truth-values}); \\ (\alpha \rightarrow \omega \rightarrow \mathbf{H}) \in \mathbb{T}_{\mathcal{F}}^{\omega} & (\text{unrestricted concepts}); \\ (\omega \rightarrow \mathbf{H}) \in \mathbb{T}_{\mathcal{F}}^{\omega} & (\text{restricted concepts}); \\ \omega \in \mathbb{T}_{\mathcal{F}}^{\omega} & (\text{worlds}). \end{cases}$$

A derivation of (10) would now look like the following:

$$\frac{\frac{\Gamma \vdash J : \alpha \rightarrow \omega \rightarrow \mathbf{H} \quad \Gamma \vdash F : \alpha}{\Gamma \vdash JF : \omega \rightarrow \mathbf{H}} \quad \Gamma \vdash W : \omega}{\Gamma \vdash JFW : \mathbf{H}} \quad (\rightarrow e) \quad (11)$$

The purpose of W is to pick out one (or several, if there are no ambiguities) valuating function(s). So, in **Example 1**, W picks out $v_{\mathcal{L}}$, in **Example 3**, W picks out $\{v_{\mathcal{L}}, v_0\}$, and in **Example 2**, W either picks out $v_{\mathcal{L}}$ or v_0 but not both. A benefit of $W : \omega$ is that it can pick out the meta-language (v_0), meta-meta-language (v_1), meta-meta-meta-language (v_2), etc. (v_n). Furthermore, W can also pick out combinations of languages ($\{v_{\mathcal{L}}, v_0, v_2, v_n\}$), provided the truth value stays consistent.

Remark 6. A caveat of this syntactic addition is that W tends to be implicit and thus, the concept horse paradox merely shifts (now we have an ambiguity of W) and does not

completely disappear. A rather elegant fix is letting W pick out $v_{\mathcal{L}}$ by default. In such a case, the paradox would dissolve and we would still preserve the truth-value in sentences like **Example 3**.