

Frege: Logical objects by abstraction and their criteria of identity

Matthias Schirn (University of Munich, Munich Center of Mathematical Philosophy)

Abstraction à la Frege

A schema for a Fregean abstraction principle can be stated as follows:

$$(S) Q(\ulcorner \urcorner) = Q(\ulcorner \urcorner) \ulcorner R_{eq}(\ulcorner \urcorner, \ulcorner \urcorner).$$

Here “ Q ” is a singular term-forming operator, $\ulcorner \urcorner$ and $\ulcorner \urcorner$ are free variables of the appropriate type, ranging over the members of a given domain, and “ R_{eq} ” is the sign for an equivalence relation holding between the values of $\ulcorner \urcorner$ and $\ulcorner \urcorner$. In *The Foundations*, the paradigms of Fregean abstraction are:

$$(1) D(a) = D(b) : \ulcorner a \parallel b.$$

The direction of line a is identical with the direction of line b if and only if line a is parallel to line b .

$$(2) N_x F(x) = N_x G(x) \ulcorner Eq_x(F(x), G(x)).$$

The number of F s is identical with the number G s if and only if F and G are equinumerous, i.e., according to Frege's definition of equinumerosity in §72 of *The Foundations*: if and only if there is a relation R that correlates one-to-one the F s and the G s. As usual, I refer to (2) as “Hume's Principle”.

In *Basic Laws*, Frege's paradigm is

$$(3) \ulcorner f(\ulcorner \urcorner) = \ulcorner g(\ulcorner \urcorner) \ulcorner \ulcorner x(f(x) \ulcorner g(x)).$$

The course-of-values of the function f is identical with the course-of-values of the function g if and only if f and g are coextensive.

Principle (3) is the famous-infamous Axiom V of Frege's logical system in *Basic Laws*, the exact structural analogue of (2).

As to Frege's idea towards the end of *The Foundations* that extensions of concepts could be dispensed with in pursuit of the logicist programme, there are just two options that he might have had in mind: (a) to identify the cardinals with objects other than extensions of concepts

or (b) to resume the tentative contextual definition of the cardinality operator (see below) and make an additional stipulation without invoking extensions of concepts. I analyze (a) and (b).

Hume's Principle and higher-order abstraction in The Foundations

In order to prove Hume's Principle, Frege has to show, according to the final explicit definition of the cardinality operator (see below), that (1) $Eq_x(F(x),G(x)) \sqsubseteq \# \sqsubseteq (Eq_x(\sqsubseteq(x),F(x))) = \# \sqsubseteq (Eq_x(\sqsubseteq(x),G(x)))$; that is to say, he has to prove that "under this hypothesis" the following two sentences hold generally: (2) $Eq_x(H(x),F(x)) \sqsubseteq Eq_x(H(x),G(x))$ and (3) $Eq_x(H(x),G(x)) \sqsubseteq Eq_x(H(x),F(x))$. Thus, Frege is converting here the statement that the extensions of two specific second-level concepts are identical into the statement that these concepts are coextensive: $\# \sqsubseteq (Eq_x(\sqsubseteq(x),F(x))) = \# \sqsubseteq (Eq_x(\sqsubseteq(x),G(x))) \sqsubseteq \sqsubseteq \sqsubseteq (Eq_x(\sqsubseteq(x),F(x)) \sqsubseteq Eq_x(\sqsubseteq(x),G(x)))$. We might thus presume that in introducing extensions of second-level concepts he had in mind the following third-order abstraction principle (modelled on Axiom V):

$$\#f(\sqsubseteq \sqsubseteq f(\sqsubseteq)) = \#f(\sqsubseteq \sqsubseteq f(\sqsubseteq)) \sqsubseteq \sqsubseteq f(\sqsubseteq \sqsubseteq f(\sqsubseteq)) \sqsubseteq \sqsubseteq \sqsubseteq f(\sqsubseteq).$$

Attempted inductive definition of the natural numbers (the finite cardinals): the number 0 belongs to the concept F, etc. (§55)

$$(I) \quad N_x^0 F(x) := \sqsubseteq x \neg F(x);$$

$$(II) \quad N_x^1 F(x) := \sqsubseteq \sqsubseteq x \neg F(x) \sqsubseteq \sqsubseteq x \sqsubseteq y (F(x) \sqsubseteq F(y) \sqsubseteq x = y);$$

$$(III) \quad N_x^{n+1} F(x) := \sqsubseteq x (F(x) \sqsubseteq N_y^n (F(y) \sqsubseteq y \neq x)).$$

The tentative contextual definition of the cardinality operator "the number which belongs to the concept \sqsubseteq " in terms of Hume's Principle (cf. §65)

$$N_x F(x) = N_x G(x) := Eq_x(F(x),G(x)).$$

The final explicit definition of the cardinality operator (§68):

$$N_x F(x) := \# \sqsubseteq (Eq_x(\sqsubseteq(x),F(x))).$$

The number which belongs to the concept F is the extension of the concept *equinumerous with the concept F*.

Cardinal arithmetic in The Foundations, some further definitions:

n is a cardinal number ($N(n) := \# \{x \mid N_x(x) = n\}$)

$0 := N_x(x \neq x)$;

$1 := N_x(x = 0)$;

$\omega_0 := N_x FN(x)$

n is a finite cardinal number ($FN(n) := n$ belongs to the natural series of numbers beginning with 0.

Hume's Principle, taken jointly with the explicit definition of the cardinality operator and the definitions of "0", "1", etc. (note that, according to Frege, each finite cardinal number except 0 can be defined as the number belonging to the concept under which just its predecessors fall) and " N_x " provides a means of determining the truth-value of any equation of the following six types (exceptions aside): (1) $N_x F(x) = N_x G(x)$; (2) $N_x F(x) = \# \{x \mid Eq_x(x, G(x))\}$; (3) $\# \{x \mid Eq_x(x, F(x))\} = \# \{x \mid Eq_x(x, G(x))\}$; (4) $N_x F(x) = n$; (5) $n = \# \{x \mid Eq_x(x, F(x))\}$; (6) $n = m$. An equation of type (4), (5), or (6) can be reduced to one of type (1), (2), or (3), and the truth-conditions of the latter three are determined via the right-hand side of Hume's Principle, or, spelled out more fully, via the one-to-one correlation of the objects falling under $F(x)$ with those falling under $G(x)$.

"Higher" numbers by abstraction in The Foundations

In *Frege 1884*, §104, Frege deals briefly with fractions, irrational and complex numbers. Just as in the case of cardinal numbers, here, too, he write

everything will in the end depend on the search for a judgeable content which can be transformed into an equation, whose sides are just the new numbers. In other words, we must fix the sense of a recognition-judgement for such numbers. In doing so, we must bear in mind the doubts that we discussed (in §§63-68) with respect to such a transformation. If we follow the same procedure as we did there, then the new numbers will be given to us as extensions of concepts.

In a first step, Frege has to contrive a suitable equivalence relation R_{eq} for the case of, let us say, the real numbers, which can be defined in purely logical vocabulary. In a second step, the real numbers are tentatively introduced by transforming $R_{eq}(\square, \square)$ into an identity of real

numbers $\lambda(\lambda) = \lambda(\lambda)$ and by presenting this transformation as a contextual definition of the λ operator. For the sake of convenience, I refer to the hypothetical abstraction principle for the reals “ $\lambda(\lambda) = \lambda(\lambda) \square \text{Req}(\lambda, \lambda)$ ” as “AR”.

The domain of the first-order variables in Basic Laws and the referential indeterminacy of course-of-values terms

Assumption “(A)” (to be discussed at length): ***A number of remarks Frege makes in The Basic Laws suggest that he takes the first-order domain to be unrestricted or all-inclusive.***

The permutation argument in *Basic Laws*, §10 (however, in the talk I shall only mention it):

Suppose h is a one-to-one function of all objects of the first-order domain of Frege’s logical system. Then the criterion of identity that applies to the references of course-of-values terms, namely the coextensiveness of two functions $\lambda(\lambda)$ and $\lambda(\lambda)$, also applies to the references of function-value names of the form “ $h(\lambda(\lambda))$ ”. On this assumption, the following equation is true:

$$(1) \quad (\lambda(\lambda) = \lambda(\lambda)) \square (h(\lambda(\lambda)) = h(\lambda(\lambda))).$$

From (1) and Axiom V follows

$$(2) \quad (h(\lambda(\lambda)) = h(\lambda(\lambda))) \square \square x(\lambda(x) \square \lambda(x)).$$

Consequently, so Frege argues, Axiom V fails to fix completely the reference of a course-of-values term “ $\lambda(\lambda)$ ”, at least if there is a bijection h such that for some course-of-values, say $\lambda(\lambda)$, $\lambda(\lambda) \neq h(\lambda(\lambda))$. In more modern terms, Frege’s argument can be presented as follows:

(P1) Suppose λ is the intended assignment of objects to course-of-values terms satisfying Axiom V. Let h be a non-trivial permutation (of all objects), and consider the assignment λ' of objects to course-of-values terms which is related to λ as follows: If λ is assigned by λ to a given course-of-values term and $\lambda = h(\lambda)$, then λ is assigned by λ' to that course-of-values term. It follows that λ' is an assignment of objects to course-of-values terms distinct from λ , but such that it satisfies Axiom V if λ does.

A special variant of (P1):

(P2) As in (P1), let λ be an assignment of objects to course-of-values terms satisfying Axiom V. Let $f(\lambda)$ and $g(\lambda)$ be two particular, extensionally non-equivalent functions. Let λ assign a to “ $\lambda f(\lambda)$ ” and b to “ $\lambda g(\lambda)$ ”. Let h be a function such that

- (i) $h(a)$ is the True,
- (ii) $h(\text{the True})$ is a ,
- (iii) $h(b)$ is the False,
- (iv) $h(\text{the False})$ is b , and,
- (v) for every argument x distinct from these, $h(x) = x$.

Finally, let σ' be an assignment of objects to course-of-values terms related to σ as in (P1), with respect to the particular permutation h just specified. Then, as in (P1), σ' will satisfy Axiom V.

In §10, the True and the False are identified with their unit classes by invoking the permutation argument. The True is identified with $\neg(\neg\top)$ — the horizontal function $\neg\top$ is elucidated as a concept under which only the True falls; it is obviously coextensive with $\top = (\top = \top)$ — and the False with $\neg(\top = \top \wedge x(x = x))$ (I use modern notation for Frege's symbols for negation and the universal quantifier).

In one place of my talk, I consider a hypothetical (third primitive) monadic second-level function in Frege's logical system: $\hat{\lambda}x.\lambda y.(x = y)$. I refer to it as 'Russell function'.

One example in support of my claim that not only in The Foundations but also in Basic Laws Frege takes the first-order domain to be all-inclusive, at least when he comes to elucidate and define first-level functions. (For reasons of time, I shall only mention the example in my talk).

The instructive example I have in mind is Frege's definition of the "membership-function" (the relation of an object falling within the extension of a concept) $\lambda x.\lambda y.(x = y)$ in §34.

$$a \in u := \forall y[\lambda x.\lambda y.(x = y)(a) = y].$$

(According to Frege's stipulations concerning the (definite) descriptive function $\lambda x.\lambda y.(x = y)$ in §11, it holds: If the function $\lambda x.\lambda y.(x = y)$ is a concept under which exactly one object falls, then " $\hat{\lambda}x.\lambda y.(x = y)$ " refers to that very object. In the remaining three cases — if more than one object falls under $\lambda x.\lambda y.(x = y)$; if no object falls under $\lambda x.\lambda y.(x = y)$; if $\lambda x.\lambda y.(x = y)$ is not a concept — " $\hat{\lambda}x.\lambda y.(x = y)$ " has the same reference as " $\lambda x.\lambda y.(x = y)$ ".)

Before setting up the definition of $\lambda x.\lambda y.(x = y)$ which involves second-order quantification, Frege stresses that $\lambda x.\lambda y.(x = y)$ must be explained (defined) for all possible objects as arguments. I take the phrase "for all possible objects as arguments" to mean: for all objects whatever as arguments. It cannot plausibly mean: for all courses-of-values as arguments. Here are my

reasons. After having defined $\lambda x. \lambda y. x$ Frege gives a number of explications and closes §34 by summarizing:

two cases must be distinguished if the value of the function $\lambda x. \lambda y. x$ is to be determined. If the λ -argument is a course-of-values, then the value of the function $\lambda x. \lambda y. x$ is the value of that function whose course-of-values is the λ -argument for the λ -argument as argument. If, on the other hand, the λ -argument is not a course-of-values, then the value of the function $\lambda x. \lambda y. x$ is $\uparrow(\neg x = \perp)$ for every λ -argument [that is, the extension of an empty concept].

Now if Frege had intended to define $\lambda x. \lambda y. x$ only for courses-of-values and the two truth-values qua special courses-of-values as arguments, he could and should have confined himself to stating what the value of $\lambda x. \lambda y. x$ is, if the λ -argument is a course-of-values. Yet he does not proceed in this way. “If, on the other hand, the λ -argument is not a course-of-values” cannot mean: “if the λ -argument is a truth-value”. In the light of the identification of the True and the False with their unit classes, this would not make sense, because the value of $\lambda x. \lambda y. x$ for the two truth-values (their unit classes) as λ -argument is already determined via the stipulation of the first case: it is the value of that function whose course-of-values is the λ -argument for the λ -argument as argument. Consequently, nor can it mean “if the λ -argument is a truth-value or any other object distinct from both the True and the False and any course-of-values”. It can only mean: if the λ -argument is any object distinct from any course-of-values.

If for Frege a sound elucidation of “ $\uparrow \lambda x. \lambda y. x$ ” had been feasible, that is, one which did not rest on a presupposed acquaintance with courses-of-values, then he could have defined straight away the predicate “ a is a course-of-values” (“ $CV(a)$ ”), modelled on his definition of “ n is a cardinal number” in *Frege 1884*, §72:

$$CV(a) := \lambda x. \lambda y. \lambda z. (x \dot{\cup} y \dot{\cup} z = a).$$

Frege’s proof of referentiality in Basic Laws, §31 — the case of the course-of-values operator “ $\uparrow \lambda x. \lambda y. x$ ”

In *Basic Laws*, §29, Frege states five criteria of referentiality for concept-script names of five different categories. With respect to “ $\uparrow \lambda x. \lambda y. x$ ”, it suffices to mention the first criterion (A) for monadic first-level function-names, the second (B) for proper names and the fourth (D) for monadic second-level function-names. For the sake of perspicuity, I present these criteria in a formalized version and apply them directly to “ $\uparrow \lambda x. \lambda y. x$ ”. I use “ λ ” as a semantic predicate

