



Swansea University
Prifysgol Abertawe



Cronfa - Swansea University Open Access Repository

This is an author produced version of a paper published in :
Logic Colloquium 2004

Cronfa URL for this paper:
<http://cronfa.swan.ac.uk/Record/cronfa13>

Book chapter :

Setzer, A. (2009). *Universes in type theory I - inaccessible and Mahlo*. Logic Colloquium 2004, (pp. 123-156).
Cambridge: Cambridge University Press.

This article is brought to you by Swansea University. Any person downloading material is agreeing to abide by the terms of the repository licence. Authors are personally responsible for adhering to publisher restrictions or conditions. When uploading content they are required to comply with their publisher agreement and the SHERPA RoMEO database to judge whether or not it is copyright safe to add this version of the paper to this repository.
<http://www.swansea.ac.uk/iss/researchsupport/cronfa-support/>

Universes in Type Theory Part I – Inaccessibles and Mahlo

Anton Setzer*

July 8, 2005

Abstract

We give an overview over universes in Martin-Löf type theory and consider the following universe constructions: a simple universe, E. Palmgren’s super universe and the Mahlo universe. We then introduce models for these theories in extensions of Kripke-Platek set theory having the same proof theoretic strength. The extensions of Kripke-Platek set theory used formalise the existence of a recursively inaccessible ordinal, a recursively hyper-inaccessible ordinal, and a recursively Mahlo ordinal. Using these models we determine upper bounds for the proof theoretic strength of the theories in questions. In case of simple universes and the Mahlo universe, these bounds have been shown by the author to be sharp. This article is an overview over the main techniques in developing these models, full details will be presented in a series of future articles.

1 Introduction

This article presents some results of a research program with the goal of determining as strong as possible predicatively justified extensions of Martin-Löf type theory (MLTT) and to determine their precise proof theoretic strength. We see three main reasons for following this research program:

1. The goal is to develop type theoretic analogues of the theories analysed in proof theory. This way we hope to make the rather abstract and technically difficult results from proof theory more accessible to the general audience, and we hope as well to give some computational meaning to those results. This will be particularly important in case of the Π_3 -reflecting universe (to be presented in the followup article [32]) which was developed from Rathjen’s ordinal notation system for $KP + (\Pi_3 - \text{refl})$ (Kripke-Platek set theory extended by the principle of Π_3 -reflection).

*Department of Computer Science, University of Wales Swansea, Singleton Park, Swansea SA2 8PP, UK, Email: a.g.setzer@swan.ac.uk, <http://www.cs.swan.ac.uk/~csetzer/>, Tel: +44 1792 513368, Fax: +44 1792 295651. Supported by Nuffield Foundation, grant ref. NAL/00303/G and EPSRC grant GR/S30450/01.

2. One can consider predicatively justified extensions of MLTT as safe theories, namely theories with a philosophical argument which justifies the validity of everything shown in those theories. This philosophical argument is given by Martin-Löf's meaning explanations. The proofs of the lower bounds show that these extensions of MLTT show the consistency of corresponding extensions of Kripke-Platek set theory, or more precisely of any approximation of it. If one accepts meaning explanations as a philosophical argument for the validity of statements shown in MLTT, one obtains in this way a consistency proof for strong extensions of Kripke-Platek set theory, which are strong enough to prove a large portion of mathematical theories (see the results in reverse mathematics, e.g. Simpson's book [33], which show that rarely more than the strength of $(\Pi_1^1 - CA)_0$ is needed; one exception seems to be the graph minor theorem). These results are in line with a revised Hilbert's program. Hilbert's original program was to prove the consistency of axiom systems for formalising mathematical proofs by finitary arguments. By Gödel's incompleteness theorem we know that this program cannot be carried out. MLTT can be considered as one replacement for finitary methods in a revised Hilbert's program, and our research demonstrates that MLTT can well be used for this purpose.

We should note however that meaning explanations haven't been worked out yet in the case of the Mahlo universe – the author himself doesn't have a sufficient background in philosophy to fill in this gap.

3. We hope as well that the new sets developed can be used as data structures in general computing. This hope has been fulfilled in case of the Mahlo universe. The data type of inductive recursive definitions used in the closed formalisation of inductive-recursive definitions developed by P. Dybjer and the author (see [10, 12, 11, 13]) has a similar character as the Mahlo universe. Variants of this data type have been used in [9, 7] in the area of generic or generative programming with the goal of writing programs which not only use data but as well the buildup of data structures given to them in order to compute new data structures and data. We hope that the Π_3 -reflecting universe, which will be developed in the followup article [32], will give rise to similar new data structures.

We should note that whether the Mahlo universe and extensions of it presented in the follow up article [32] are actually predicatively justifiable is still a matter of a debate. However, even if there is a debate on this, it seems at least at the moment to be unlikely that a better construction can be found which avoids the problem in this case. For instance Erik Palmgren's higher order universes [20] look at first sight more directly predicatively justifiable – but a closer look (it took the author a long time to actually discover this) reveals that in Palmgren's construction Set has a similar character as a Mahlo universe: The introduction rule allows to introduce new elements of Set assuming an arbitrary function from families of Set into families of Set (the new element will be a universe closed under this function). This is essentially the same as the introduction rule

for the Mahlo universe.

Follow up articles. The original goal when writing this article was to present the Π_3 -reflecting universe. It turned out that much more space is needed to present it, and therefore the presentation of the Π_3 -reflecting universe will be given in a planned followup article [32]. Even so the current article is rather long, we won't have room to present the models in full detail. We plan to write a series of articles, in which the models will be worked out in full detail. These future articles will then set up a better infrastructure for writing future articles on model constructions for MLTT. So the current article serves as an overview article, which presents the basic techniques for developing models of type theory suitable for determining upper bounds for the proof theoretic strength. It might be that this article is more accessible than the more detailed technical articles to follow.

Content. The structure of this article is as follows: In Sect. 2 we briefly develop the small and large logical framework and the basic set constructions (i.e. the sets N_n , N and the sets formed using $+$, Σ , Π , W , Id). We then briefly introduce Kripke-Platek set theory. Then we develop the basic principles for defining models of type theory for the basic set constructions and the small logical framework. In Sect. 3 we introduce the rules for universes and the theory ML_WU . Furthermore we determine the principles for introducing nested universes. We will see that there are two ways of defining subuniverses: recursive and inductive subuniverses. We then develop the main concepts for introducing models for universes and develop in particular a model of ML_WU . In Sect. 4 we introduce Erik Palmgren's super universe and develop a model of it in a corresponding extension of Kripke-Platek set theory. Finally, in Sect. 5 we introduce the Mahlo universe and determine a model of it. For all three universe constructions mentioned we obtain proof theoretic upper bounds for the strength of the theories in question.

Related results. The first model of MLTT developed for proof theoretic purposes was the model of MLTT with one universe but no W -type, developed by P. Aczel ([1]). Rathjen and Griffor have analysed the strength of MLTT with finitely many universes, W -type, and of several variants of it in [14]. Erik Palmgren has introduced the super universe and later higher order universes (the latter ones are conjectured to reach the strength of KPM), which are best presented in [20]. There is a rich literature on PER models of MLTT, an overview is given in M. Hofmann's article [15]. Note that the emphasis of this article is in developing models in order to determine upper bounds for the proof-theoretic strength of those theories, which means as well to develop models using minimal strength in the Meta-theory. In this research programme we are following the steps of ordinal theoretic proof theory, in which the major recent steps were done by M. Rathjen, who analysed the theories KPM ([22]), $\text{KP} + (\Pi_3 - \text{refl})$ [23] and $\Pi_2^1 - \text{CA} + \text{BI}$ ([24]). Similar steps have been taken by T. Arai ([3, 2, 4, 5];

there are as well preprints of Arai covering Π_n -reflection, Σ_1 -stability and Π_1 -collection).

This research benefitted very much from intensive discussions with T. Arai, U. Berger, P. Martin-Löf, E. Palmgren, and M. Rathjen.

2 Basic Martin-Löf Type Theory

In this section we will briefly introduce the version of Martin-Löf type theory (MLTT) used in this article. We will require some knowledge about MLTT. The reader not familiar with it might refer to [17, 18, 19, 21] or Chapter 11 of [34] (the latter deviates in its description of intensional type theory from the standard versions of MLTT).

2.1 The Small Logical Framework

The role of the logical framework in this article. In this article we will, when formally introducing theories and analysing them, not make use of the logical framework. There are two reasons for it:

On one hand, the author has at the moment conceptual problems with the logical framework. The problem manifests itself particularly when considering meaning explanations. Meaning explanations for the logical framework don't seem to have been worked out fully at present, whereas the concept of meaning explanations for type theory without the logical framework seems to be well understood. One should note however that Per Martin-Löf seems to have a clear understanding of meaning explanations for the logical framework, and has given talks on this topic.

On the other hand, a full treatment of the logical framework causes at the moment still problems when determining proof theoretic bounds. Our techniques for modelling type theories at present don't allow to model it directly without using more strength than is actually needed. Therefore, with our techniques we won't be able to obtain precise proof theoretic bounds. In order to avoid this, it seems to be necessary to first eliminate the use of the logical framework, and then to use the techniques used in this article. Martin Hofmann has shown in Sect. 4.3 of [15] that type theory with the logical framework is conservative over type theory without it, which would give the desired reduction. However, we haven't had time yet to study this result in detail yet, in order to make sure that it can be applied to our setting.

For these two reasons all theories presented in the following will only make use of the restriction of the logical framework to set, which we call the small logical framework. However, in order to explain the heuristics of our constructions, it is sometimes useful to make use of the full logical framework. We will do so for heuristic purposes only – when the formal theories are introduced, we will not make use of it.

The small logical framework. The small logical framework contains the dependent function set and product as in the the full logical framework, but limited to sets. So we have, under the assumptions $A : \text{Set}$, $x : A \Rightarrow B : \text{Set}$ the following set constructions:

- The dependent function set $(x : A) \rightarrow B : \text{Set}$.
 - The introduction rule expresses that we can form $(x : A)t : (x : A) \rightarrow B$ provided $x : A \Rightarrow t : B$.
 - The elimination rule expresses that we can apply $f : (x : A) \rightarrow B$ to $a : A$ and obtain $f(a) : B[x := a]$.
 - We have as equality rule β -equality, i.e. $((x : A)t)s = t[x := s]$.
 - Additionally we have the η -rule: if $f : (x : A) \rightarrow B$ then $f = (x : A)f(x)$.
 - Furthermore we have equality rules of the formation rule (the rule which forms $(x : A) \rightarrow B : \text{Set}$, provided $A : \text{Set}$ and $x : A \Rightarrow B : \text{Set}$), the introduction rule and the elimination rule:
 - * The equality version of the formation rule expresses that if $A = A' : \text{Set}$, $x : A \Rightarrow B = B' : \text{Set}$, then $(x : A) \rightarrow B = (x : A') \rightarrow B' : \text{Set}$.
 - * The equality version of the introduction rule expresses that if $x : A \Rightarrow t = t' : B$ then $(x : A)t = (x : A)t'$ (that's the ξ -rule).
 - * The equality version of the elimination rule expresses that if $s = s' : (x : A) \rightarrow B$ and $t = t'$ then $s(t) = s'(t')$.
 - * The principle of forming the equality versions of the formation, introduction and elimination rules from the standard (non-equality) rules is a straightforward principle. Therefore the convention is that when we introduce rules in the following, we silently introduce as well the equality versions as well. An exception are the equality rules, which don't have an equality version.
 - * α -equivalent terms are considered to be identical, therefore there is no explicit rule for α -equality. This extends to judgements as well, i.e. $x : A, y : B \Rightarrow C : \text{Set}$ and $u : A, v : B[x := u] \Rightarrow C[x := u, y := v] : \text{Set}$ are considered as the same judgement.
 - We mention the following abbreviations:
 - * We write $f(a, b, c)$ instead of $f(a)(b)(c)$, similarly for longer applications.
 - * We usually omit the type A in $(x : A)t$ and write $(x)t$ instead.
 - * $(x, y)t := (x)(y)t$, similarly for longer expressions.
 - * We write $(x, y : A) \rightarrow B$ for $(x : A) \rightarrow (x : B) \rightarrow B$, $(x : A, y : B) \rightarrow C$ for $(x : A) \rightarrow (y : B) \rightarrow C$. Similar abbreviations are to be understood in the same way.
 - * If $f : A \rightarrow B$ and $g : B \rightarrow C$ then $g \circ f := (x)g(f(x))$.

- The dependent product $(x : A) \times B : \text{Set}$.
 - The introduction rule allows to form $\langle a, b \rangle$ for $a : A$ and $b : B[x := a]$.
 - The elimination rule forms the projections of an element $c : (x : A) \times B : \text{Set}$, written as $\pi_0(c) : A$ and $\pi_1(c) : B[x := \pi_0(c)]$.
 - The equality rule expresses that $\pi_0(\langle a, b \rangle) = a$ and $\pi_1(\langle a, b \rangle) = b$.
 - We have as well the η -rule associated with it, so if $c : (x : A) \times B$ then $c = \langle \pi_0(c), \pi_1(c) \rangle$.

We introduce here as well $A \rightarrow B := (x : A) \rightarrow B$ and $A \times B := (x : A) \times B$ for some fresh variable x .

Definition 2.1. *By the rules of the small logical framework we mean the structural rules of type theory and the rules for the dependent function set and the dependent product restricted to Set.*

2.2 The Full Logical Framework

The (full) logical framework (which is never part of the official type theories in this article), is obtained in the following way: Apart from the type sets, which is the highest type in the version excluding the logical framework, one has as well a type of large sets Type . Type contains Set and every element of Set . So we have the rules

$$\text{Set} : \text{Type} \quad \frac{A : \text{Set}}{A : \text{Type}} \quad \frac{A = B : \text{Set}}{A = B : \text{Type}}$$

Apart from Set , Type will as well be closed under the dependent function type and the dependent product. So if $A : \text{Type}$ and $x : A \Rightarrow B : \text{Type}$, we have that

- $(x : A) \rightarrow B : \text{Type}$,
- $(x : A) \times B : \text{Type}$,

with essentially the same rules as the small logical framework, except that one refers to Type instead of Set .

2.3 The Basic Set Constructions

The basic set constructions are the following sets (or more precisely principles for forming sets) and their corresponding rules:

- The finite sets $N_n : \text{Set}$, where $n \in \{0, 1, \dots, \}$; note that here n is not an internal natural number inside type theory, but exists on the Meta-level. The introduction rules are $A_i^n : N_n$ for $i = 0, \dots, n - 1$ (where i are again numbers on the Meta-level).
- The set N of natural numbers. The introduction rules are $0 : N$ and $S : N \rightarrow N$.

- The disjoint union $A + B$ of two sets A, B , with introduction rules $\text{inl} : A \rightarrow (A + B)$ and $\text{inr} : B \rightarrow (A + B)$.
- For $A : \text{Set}$ and $x : A \Rightarrow B : \text{Set}$ we have the following sets:
 - The Π -set $\Pi x : A.B$. The introduction rule is $\lambda : ((x : A) \rightarrow B) \rightarrow \Pi x : A.B$. One writes $\lambda x : A.t$ for $\lambda((x : A)t)$.
 - The Σ -set $\Sigma x : A.B$. The introduction rule is $\text{p} : (a : A) \rightarrow (b : B[x := a]) \rightarrow \Sigma x : A.B$. (The differences between $\Pi x : A.B$, $\Sigma x : A.B$, and $(x : A) \rightarrow B$, $(x : A) \times B$ will be explained below).
 - The W set $\text{W}x : A.B$, which is the set of well-founded trees with branching degrees $B[x := a]$ for $a : A$. The introduction rule is $\text{sup} : (a : A, b : B[x := a]) \rightarrow \text{W}x : A.B \rightarrow \text{W}x : A.B$.
- The intensional identity set $\text{Id}(A, a, b)$ for $A : \text{Set}$, $a : A$, $b : A$. The introduction rule is $\text{refl}_A : (a : A) \rightarrow \text{Id}(A, a, a)$.
- For all of the above set constructions, the elimination rules express that the above sets are the least sets introduced by these constructors. This will be for instance in case of N primitive recursion into arbitrary sets (which might depend on the element $n : \text{N}$ we are eliminating), in case of W induction over those trees, and in case of $A + B$ case distinction on whether $ab : A + B$ is of the form $\text{inl}(a)$ or $\text{inr}(b)$.
- Furthermore, we have the standard equality rules for the above mentioned sets.

The main difference between the dependent function set and product and the sets $\Pi x : A.B$, $\Sigma x.B$ is that we have the η -rule for the constructions of the logical framework, but not for Π and Σ . The conceptual reason for this becomes clear when considering inductive-recursive definitions. All the above constructions are inductive-recursive (in case of $\text{Id}(A, a, b)$ general indexed inductive-recursive definitions) as introduced originally due to P. Dybjer with a formalisation using finitely many rules by P. Dybjer and the author; see [10, 8, 12, 11, 13]). Indexed inductive-recursive definitions allow to introduce all sets in MLTT by determining their introduction rules (there are of course restrictions on which introduction rules are allowed). The elimination and equality rules are then derived automatically. There is no η -rule involved in this schema, which would be unnatural in general. For instance, it does not make much sense to consider an η -rule for N. Therefore it is natural to exclude the η -rule from Π and Σ , and to have separate logical framework versions which contain the η -rule. In this article, the difference between the logical framework set constructions and the Π - and Σ -set won't play a big rôle.

Our models will admit as well the addition of an extensional equality set to the type theory. Since lower bounds will be obtained using intensional equality

only, it will follow that the proof theoretic strength of the type theories under consideration with intensional equality and with extensional equality coincides.

Definition 2.2. *By the basic set constructions in type theory we mean the above mentioned constructions $N_n, N, +, \Sigma, \Pi, W, \text{Id}$ for forming sets and the corresponding formation/introduction/elimination/equality rules.*

2.4 Kripke-Platek Set Theory

In this article, we will develop models of the type theories considered in versions of Kripke-Platek set theory. This will allow us to determine upper bounds for the theories in question.

Kripke-Platek set theory (KP) is a weak version of set theory, based on classical logic. The “bible” of KP is the book by Barwise [6]. KP can be used in order to develop most concepts in generalised recursion theory, as demonstrated in that book. KP has been pioneered by Jäger [16] as a reference theory for proof theoretic studies. For many theories, there exists a version of KP of equal strength, and often one can determine upper bounds for theories by determining upper bounds for that variant for KP and then by modelling the original theory in that variant. This is as well the approach taken in this article.

We don’t have room to introduce KP in full (we highly recommend the reader not familiar with it to study the first chapters of [6]). We briefly repeat here its axioms (see p. 10 of [6]; the theory presented in [6] adds to KP as well urelemente, and is therefore called KPU for KP plus urelemente; the versions of KP used in this article don’t have urelemente):

- (*Extensionality*) $\forall x, y. (\forall z. z \in x \leftrightarrow z \in y) \rightarrow x = y$
- (*Foundation*) $\forall \vec{z}. (\forall x. (\forall y \in x. \varphi(y, \vec{z})) \rightarrow \varphi(x, \vec{z})) \rightarrow \forall x. \varphi(x, \vec{z})$
where $y \notin \text{FV}(\varphi(x, \vec{z}))$
- (*Pair*) $\forall x, y. \exists z. x \in z \wedge y \in z$
- (*Union*) $\forall x. \exists y. \forall z \in x. \forall u \in z. u \in y$
- (Δ_0 -*Separation*) $\forall \vec{x}, y. \exists z. \forall u. (u \in z \leftrightarrow (u \in y \wedge \varphi(x, \vec{x})))$
where $\varphi(x, \vec{x})$ is Δ_0 , $z \notin \text{FV}(\varphi(x, \vec{x}))$
- (Δ_0 -*Collection*) $\forall \vec{x}, y. (\forall z \in y. \exists u. \varphi(z, u, \vec{x})) \rightarrow \exists v. \forall z \in y. \exists u \in v. \varphi(z, u, \vec{x})$
where $\varphi(z, u, \vec{x})$ is Δ_0 , $v \notin \text{FV}(\varphi(z, u, \vec{x}))$

The versions of KP used in this article will always be augmented by the axiom of infinity

- (*Infinity*) $\exists x. \emptyset \in x \wedge \forall y \in x. y \cup \{y\} \in x$,

where $\emptyset \in x$ and $y \cup \{y\} \in x$ are to be understood in the usual way. Let $\text{KP}\omega$ be the theory KP plus the axiom of infinity. When forming extensions of $\text{KP}\omega$, one adds a predicate $\text{Ad}(x)$ for “ x is an admissible containing ω ”, following the approach taken by Jäger (e.g. [16]). Here an admissible is a transitive inner model of KP, and by admissible $> \omega$ we mean an admissible containing ω , i.e.

a transitive inner model of $KP\omega$. So one has the following additional axioms ($\text{trans}(x)$ expresses that x transitive):

- (Ad.1) $\forall x. \text{Ad}(x) \rightarrow \text{trans}(x)$
- (Ad.2) $\forall x, y. \text{Ad}(x) \wedge \text{Ad}(y) \rightarrow (x \in y \vee x = y \vee y \in x)$
- (Ad.3) $\forall x, \vec{y}. (\text{Ad}(x) \rightarrow \psi^x(\vec{y}))$
for every instance $\psi(\vec{y})$ of (*Pair*), (*Union*), (Δ_0 -*Separation*), (Δ_0 -*Collection*)

$\psi^x(\vec{y})$ is obtained by replacing all unbounded quantifiers $\forall y, \exists y$ occurring in $\psi(\vec{y})$ by $\forall y \in x, \exists y \in x$ respectively, but leaving restricted quantifiers as they are (where an unrestricted quantifier is a quantifier not of the form $\forall y \in z, \exists y \in z$). The extensions of $KP\omega$ are then obtained by adding to $KP\omega$ the above axioms for Ad and axioms expressing that there exists an admissible with sufficiently strong closure properties. The convention in this article is that whenever we introduce an extension of KP, we always add as well the infinity axiom and, if Ad is involved in any of the additional axioms, the axioms (Ad.1-3).

We introduce as well the theory KP1 which expresses that there are finitely many admissibles, but that the set theoretic universe itself is not necessarily a model of KP. The standard model of KP1 is $a := \bigcup_{n \in \omega} b_n$, where b_n is the n th admissible ($n \in \omega$). Note that a is not an admissible. Formally, KP1 is obtained by taking $KP\omega$ without (Δ_0 -*Collection*) and adding the axioms (Ad.1-3) and the following axiom (*Lim*):

$$(Lim) \quad \forall x. \exists y. \text{Ad}(y) \wedge x \in y$$

When working in extensions of $KP\omega$, we will frequently make use of the constructible hierarchy, and refer to the definition of L_α in [6]. We define as well the following:

Definition 2.3. (a) $L_{<\alpha} := \bigcup_{\beta < \alpha} L_\beta$,

$$(b) \quad L'_\alpha := L_{\omega \cdot \alpha},$$

$$(c) \quad L'_{<\alpha} := \bigcup_{\beta < \alpha} L'_\beta.$$

(d) *An admissible ordinal is an ordinal α s.t. L_α is admissible.*

We will often make use of the following essentially trivial remark:

Remark 2.4. (a) *If $b = \{x \in a \mid \varphi(x, a_1, \dots, a_n)\}$ where φ is Δ_0 and $a, a_1, \dots, a_n \in L'_\alpha \cup \{L'_\alpha\}$, then $b \in L'_{\alpha+1}$.*

$$(b) \quad \alpha \in L'_{\alpha+1}.$$

$$(c) \quad \omega \in L'_2.$$

Proof. Let Ord be the class of ordinals. (a): Assumption II.5.2 and Theorem II.6.4 of [6]. (b): One easily can show that $L_\alpha \cap \text{Ord} \in \text{Ord}$. We show the assertion by induction on α : The cases $\alpha = 0$ and $\alpha \rightarrow \alpha + 1$ follow easily or by

IH. In case α limit we have by IH $\alpha \subseteq L'_{<\alpha} \cap \text{Ord} =: \beta \in L'_{\alpha+1}$, and therefore $\alpha \in L'_{\alpha+1}$. (c): $\omega \subseteq L_{<\omega}$, $L_{<\omega} \cap \text{Ord} \subseteq \omega$, therefore $\omega := \{x \in L_{<\omega} \mid \text{Ord}(x)\} \in L_{\omega+\omega} = L'_2$. □

2.5 Models of the Basic Type Theory

As said before, we will develop models of the type theories in question in extensions of $\text{KP}\omega$. These models will extend the model of MLTT with one universe and the W-type, as worked out in the author's PhD thesis [25] (see as well [27] for a draft version of an article on it, and [26] for an article containing a basic model construction for MLTT with one Mahlo universe).

Basic principles of the model construction. The basic idea of all our models is that we define a PER model, i.e. a model in which sets are interpreted as partial equivalence relations on a set of basic objects, where in this article the set of basic objects will be a set of terms, which is a slightly simplified version of the set of terms occurring in the type theory. Let ML' the type theory for which we want to determine an upper bound for its proof theoretic strength. We will model ML' in an extension KP' of $\text{KP}\omega$. KP' will be defined in such a way that it has the same expected proof theoretic strength. Once we have modelled ML' in KP' one will be able to conclude (this step requires some techniques which can be found in [25]) that the proof theoretic strength $|\text{ML}'|$ of ML' is less than or equal the proof theoretic strength $|\text{KP}'|$ of KP' , i.e. $|\text{ML}'| \leq |\text{KP}'|$.

The soundness theorem for the model of type theory will be a Meta-theorem (otherwise we would obtain $|\text{ML}'| < |\text{KP}'|$): For Meta-every judgement $\Gamma \Rightarrow \theta$ s.t. $\text{ML}' \vdash \Gamma \Rightarrow \theta$ we have that KP' proves the correctness statement associated with $\Gamma \Rightarrow \theta$. But ML' doesn't prove that for all $\Gamma \Rightarrow \theta$ the correctness statement holds. In fact, ML' won't even express such a general correctness statement.

Let Term be the underlying set of terms occurring in ML' . We will introduce in the model a set of closed terms $\llbracket \text{Term} \rrbracket$ using essentially the same term and type constructors and destructors as in the type theory. α -equivalent terms (as well sets, context, judgements) will always be identified. However, we will, when introducing the models of the universes below occasionally throw away arguments needed for type checking purposes only (this is not really necessary, but simplifies matters).

One introduces as well a deterministic reduction relation \longrightarrow on $\llbracket \text{Term} \rrbracket$ corresponding to weak head reductions following the equality rules of the type theory. \longrightarrow^* is the reflexive and transitive closure of \longrightarrow .

Definition 2.5. (*Environments*)

- (a) An environment relative to a sequence of variables x_1, \dots, x_n is an expression $[x_1 = s_1, \dots, x_n = s_n]$, where x_i are distinct variables and $s_i \in \llbracket \text{Term} \rrbracket$.
- (b) \emptyset stands for the empty environment $\llbracket \rrbracket$.

- (c) $[x_1 = s_1, \dots, x_n = s_n], y = t := [x_1 = s_1, \dots, x_n = s_n, y = t]$, where we assume that after some renaming of variables corresponding to α -conversion in the situation used, y is different from x_i .
- (d) $\text{FV}([x_1 = s_1, \dots, x_n = s_n]) := \{x_1, \dots, x_n\}$.
- (e) $[x_1 = s_1, \dots, x_n = s_n](x_i) := s_i$.
- (f) $\text{Env}(x_1, \dots, x_n) := \{[x_1 = s_1, \dots, x_n = s_n] \mid s_1, \dots, s_n \in \llbracket \text{Term} \rrbracket\}$.

The interpretation $\llbracket t \rrbracket_\rho \in \llbracket \text{Term} \rrbracket$ for environments ρ and terms $t \in \text{Term}$ s.t. $\text{FV}(t) \subseteq \text{FV}(\rho)$ is the result of omitting in t the arguments of constructors/destructors to be thrown away (as indicated when introducing the constants) and substituting free variables x in t by $\rho(x)$.

We write $\llbracket t \rrbracket$ for $\llbracket t \rrbracket_\emptyset$.

Then one defines the interpretation $\llbracket A \rrbracket_\rho$ of sets A relative to environments ρ s.t. $\text{FV}(A) \subseteq \text{FV}(\rho)$ as a set of pairs $\langle s, t \rangle$ of closed terms. The intended meaning for $\langle s, t \rangle \in \llbracket A \rrbracket_\rho$ is that s and t are equal elements of the interpretation of A . We write $\llbracket A \rrbracket$ for $\llbracket A \rrbracket_\emptyset$.

We define some basic operations on sets $X \subseteq \llbracket \text{Term} \rrbracket^2$:

Definition 2.6. Let $X \subseteq \llbracket \text{Term} \rrbracket^2$.

- (a) The closure of X under reductions is defined as

$$\text{Clos}_\rightarrow(X) := \{\langle s, t \rangle \mid \exists \langle s', t' \rangle \in X \mid s \longrightarrow^* s' \wedge t \longrightarrow^* t'\}.$$

- (b) $\text{Flat}(X) := \{a \in \llbracket \text{Term} \rrbracket \mid \exists b \in \llbracket \text{Term} \rrbracket. \langle a, b \rangle \in X \vee \langle b, a \rangle \in X\}$.

- (c) $\text{PER}(X)$ iff X is a partial equivalence relation (i.e. symmetric and transitive) and X is closed downwards under \longrightarrow , i.e. if $s \longrightarrow^* s'$, $t \longrightarrow^* t'$ and $\langle s', t' \rangle \in X$, then $\langle s, t \rangle \in X$.

$\text{Flat}(\llbracket A \rrbracket_\rho)$ is the set of elements of the interpretation of A , and $\langle a, b \rangle \in \llbracket A \rrbracket_\rho$ means that a and b are equal elements of $\llbracket A \rrbracket_\rho$. This is meaningful if $\text{PER}(\llbracket A \rrbracket_\rho)$ holds, which we will show for all derivable sets A as expressed by the correctness statement below.

In the models for the type theories considered in this article, we can as well assign to every set A with $\text{FV}(A) \subseteq \{x_1, \dots, x_n\}$ an ordinal $\text{o}(A) \geq 3$ s.t. for all $\rho \in \text{Env}(x_1, \dots, x_n)$ we have $\llbracket A \rrbracket_\rho \in L'_{\text{o}(A)}$, independent of ρ . Note that $\llbracket \text{Term} \rrbracket$ is a Δ_0 -definable subset of ω and therefore $\llbracket \text{Term} \rrbracket \in L'_3$.

The definition of $\llbracket A \rrbracket_\rho$ and $\text{o}(A)$ for A formed by one of the standard set constructors is as follows:

- $\llbracket \mathbb{N} \rrbracket_\rho := \llbracket \mathbb{N} \rrbracket := \text{Clos}_\rightarrow(\{\langle S^n(0), S^n(0) \rangle \mid n \in \omega\})$.
 $\text{o}(\mathbb{N}) := 3$ (this is probably not the minimal ordinal one can associate with it, similarly for the later ordinals, but this bound suffices for our proof).

- $\llbracket \Sigma x : A.B \rrbracket_\rho := \llbracket \Sigma \rrbracket(\llbracket A \rrbracket_\rho, \lambda s \in \text{Flat}(\llbracket A \rrbracket_\rho) \cdot \llbracket B \rrbracket_{\rho, x=s})$,
where, if $A \subseteq \llbracket \text{Term} \rrbracket^2$, and for $a \in \text{Flat}(A)$ we have $B(a) \subseteq \llbracket \text{Term} \rrbracket^2$,
then

$$\llbracket \Sigma \rrbracket(A, B) := \text{Clos}_{\rightarrow}(\{\langle p(a, b), p(a', b') \rangle \mid \langle a, a' \rangle \in A \wedge (B(a) = B(a')) \wedge \langle b, b' \rangle \in B(a)\}) .$$

(Here $\lambda x \in A.s(x) := \{\langle x, s(x) \rangle \mid x \in A\}$).

$$o(\Sigma x : A.B) := \max\{o(A), o(B)\} + 2.$$

- Similarly one defines the semantics and the ordinals of $\Pi x : A.B$, $Wx : a.B$, $A + B$, $\text{Id}(A, a, b)$, $(x : A) \rightarrow B$, $(x : A) \times B$, making use of suitable definitions of $\llbracket \Pi \rrbracket$, $\llbracket W \rrbracket$, $\llbracket + \rrbracket$, $\llbracket \text{Id} \rrbracket$, $\llbracket \rightarrow \rrbracket$ and $\llbracket \times \rrbracket$, and incrementing the ordinals as for Σ by 2 (except for W , see below).

We introduce here as well the following notations to be used later:

- $(x \in A) \llbracket \rightarrow \rrbracket B := \llbracket \rightarrow \rrbracket(A, \lambda x \in \text{Flat}(A).B)$,
- $A \llbracket \rightarrow \rrbracket B := (x \in A) \llbracket \rightarrow \rrbracket B$ where x is fresh,
- $(x \in A) \llbracket \times \rrbracket B := \llbracket \times \rrbracket(A, \lambda x \in \text{Flat}(A).B)$,
- $A \llbracket \times \rrbracket B := (x \in A) \llbracket \times \rrbracket B$ where x is fresh.
- We can omit brackets in the above, using for $\llbracket \rightarrow \rrbracket$, $\llbracket \times \rrbracket$ the same conventions as we have for \rightarrow , \times .

Only in case of W , one defines $o(Wx : A.B) := (\max\{o(A), o(B)\})^+ + 1$, where α^+ is the next admissible ordinal above α . The reason for using $\alpha^+ + 1$ in case of $\llbracket W \rrbracket$ is that $\llbracket W \rrbracket(A, B)$ is defined by iterating an operator, which makes the one step closure under the formation of new trees with subtrees being the previously defined trees up to α^+ .

We interpret contexts as sets of environments as follows:

$$\begin{aligned} \llbracket x_1 : A_1, \dots, x_n : A_n \rrbracket := & \\ \{ \langle [x_1 = s_1, \dots, x_n = s_n], [x_1 = t_1, \dots, x_n = t_n] \rangle \mid & \\ \langle s_1, t_1 \rangle \in \llbracket A_1 \rrbracket_{\emptyset} \wedge & \\ \langle s_2, t_2 \rangle \in \llbracket A_2 \rrbracket_{[x_1=s_1]} \wedge \dots \wedge & \\ \langle s_n, t_n \rangle \in \llbracket A_n \rrbracket_{[x_1=s_1, \dots, x_{n-1}=s_{n-1}]} \} & \end{aligned}$$

Definition 2.7. *The correctness of a judgement is defined as follows by induction on the length:*

- Correct($\emptyset \Rightarrow \text{Context}$) (where \emptyset is the empty context) is the true formula.
-

$$\begin{aligned} \text{Correct}(\Gamma \Rightarrow A = B : \text{Set}) & \\ := \text{Correct}(\Gamma \Rightarrow \text{Context}) \wedge \forall \langle \rho, \rho' \rangle \in \llbracket \Gamma \rrbracket \cdot \text{PER}(\llbracket A \rrbracket_\rho) \wedge \llbracket A \rrbracket_\rho = \llbracket B \rrbracket_{\rho'} . & \end{aligned}$$

(c)

$$\begin{aligned} \text{Correct}(\Gamma, y : A \Rightarrow \text{Context}) &:= \text{Correct}(\Gamma \Rightarrow A : \text{Set}) \\ &:= \text{Correct}(\Gamma \Rightarrow A = A : \text{Set}) , \end{aligned}$$

as just defined.

(d)

$$\begin{aligned} \text{Correct}(\Gamma \Rightarrow a = b : A) \\ := \text{Correct}(\Gamma \Rightarrow A : \text{Set}) \wedge \forall \langle \rho, \rho' \rangle \in \llbracket \Gamma \rrbracket . \langle \llbracket a \rrbracket_\rho, \llbracket a \rrbracket_{\rho'} \rangle \in \llbracket A \rrbracket_\rho . \end{aligned}$$

(e) $\text{Correct}(\Gamma \Rightarrow a : A) := \text{Correct}(\Gamma \Rightarrow a = a : A)$ as just defined.

The correctness statement to be shown by induction on the derivations in ML' is, that for Meta-every judgement $\Gamma \Rightarrow \theta$ we have

$$(\text{ML}' \vdash \Gamma \Rightarrow \theta) \rightarrow \text{KP}' \vdash \text{Correct}(\Gamma \Rightarrow \theta) .$$

3 Universes and Recursively Inaccessibles

3.1 Rules for Universes

A universe (à la Tarski) in MLTT (Martin-Löf type theory) is given by a set $U : \text{Set}$ of codes for sets together with a decoding function $T : U \rightarrow \text{Set}$, which without the logical framework is given by the rule:

$$\frac{u : U}{T(u) : \text{Set}}$$

All universes in the following will be closed under the basic set constructions (i.e. $N_n, N, \Sigma, \Pi, W, +, \text{Id}$). This means that we have in case of N and Σ the rules:

$$\begin{array}{ccc} \widehat{N} : U & & T(\widehat{N}) = N \\ \\ \frac{a : U \quad b : T(a) \rightarrow U}{\widehat{\Sigma}(a, b) : U} & & T(\widehat{\Sigma}(a, b)) = \Sigma x : T(a). T(b(x)) \end{array}$$

(The codes for $N_n, \Pi, W, +, \text{Id}$ will be called $\widehat{N}_n, \widehat{\Pi}, \widehat{W}, \widehat{+}, \widehat{\text{Id}}$ respectively.)

In particular we consider the theory $\text{ML}_W U$ consisting of the small logical framework, the basic set constructions and one universe (which is by our convention closed under the standard set constructors including the W -set, but not under U itself; the subscript W stands for the closure under the W -set).

We will not add any elimination rules for universes. In case of the theory $\text{ML}_W U$ and of the super universe (as introduced in Sect. 4), they could be added without any problems, but won't add any strength to them. In case of the Mahlo universe (Sect. 5) and stronger universes, such elimination rules result in inconsistencies (see Theorem 6.1 in [20]).

Type theories with several universes. We will already here look at how to extend this approach to type theories having several universes. In this case, the codes $\widehat{N}, \widehat{\Sigma}$ etc. have special names and additional parameters so that we can distinguish them from each other. These parameters will always be omitted when moving to $\llbracket \text{Term} \rrbracket$. E.g. if we have universes V and $U_{a,b}$, we have in Term codes \widehat{N}_V and $\widehat{N}_{U,a,b}$, whereas in $\llbracket \text{Term} \rrbracket$ we only have a code \widehat{N} and have $\llbracket \widehat{N}_V \rrbracket_\rho = \llbracket \widehat{N}_{U,a,b} \rrbracket_\rho = \widehat{N}$.

By a *recursive subuniverse* (U, \widehat{T}_U) of a universe (V, T_V) we mean that we have $U : \text{Set}$, but, instead of having $T_U : U \rightarrow \text{Set}$ (or the corresponding formulation above avoiding the full logical framework), have the rule formulating the existence of a function $\widehat{T}_U : U \rightarrow V$, which recursively decodes codes by their corresponding codes in the other set. If we for $a : U$ define $T_U(a) := T_V(\widehat{T}_U(a)) : \text{Set}$, we have therefore the rules:

$$\widehat{N}_U : U \qquad \widehat{T}_U(\widehat{N}_U) = \widehat{N}_V$$

$$\frac{a : U \quad b : T_U(a) \rightarrow U}{\widehat{\Sigma}_U(a, b) : U} \quad \widehat{T}_U(\widehat{\Sigma}_U(a, b)) = \widehat{\Sigma}_V(\widehat{T}_U(a), \widehat{T}_U \circ b)$$

similarly for the other basic set constructions.

Sometimes it will not be possible to develop the subuniverses introduced in this article as recursive subuniverses. The problem occurs as soon as we have infinite chains of universes above one universe. Consider for instance a sequence of universes $(U_\alpha)_{\alpha \leq \beta}$. There is no problem, if we have a finite sequence of such universes (i.e. $\beta < \omega$): Then we can consider for $\alpha < \beta$ the universe U_α as a recursive subuniverse of $U_{\alpha+1}$, whereas the universe U_β itself is defined as a standard universe. The set corresponding to $a : U_0$ is then $T_\beta(\widehat{T}_{\beta-1}(\dots(\widehat{T}_0(a))\dots))$. If $\beta = \omega$, we would with this approach end up with a nonterminating chain

$$U_0 \xrightarrow{\widehat{T}_0} U_1 \xrightarrow{\widehat{T}_1} U_2 \xrightarrow{\widehat{T}_2} \dots$$

and we cannot determine a set corresponding to $x : U_0$.

The way around this is to follow Per Martin-Löf's approach, by forming *inductive subuniverses*. An inductive subuniverse U of a universe (V, T_V) is given by

- a set U ,
- a recursively defined decoding function $T_U : U \rightarrow \text{Set}$ and
- a *constructor* $\widehat{T}_U : U \rightarrow V$.

So $\widehat{T}_U : U \rightarrow V$ forms an additional introduction rule for V .

If we avoid the full logical framework, we get the following rules for an inductive subuniverse:

$$\begin{array}{ccc}
U : \text{Set} & & \frac{a : U}{T_U(a) : \text{Set}} \\
\frac{a : U}{\widehat{T}_U(a) : V} & & T_V(\widehat{T}_U(a)) = T_U(a) \\
\widehat{N}_U : U & & T_U(\widehat{N}_U) = N \\
\frac{a : U \quad b : T_U(a) \rightarrow U}{\widehat{\Sigma}_U(a, b) : U} & & T_U(\widehat{\Sigma}_U(a, b)) = \Sigma x : T_U(a).T_U(b(x))
\end{array}$$

Similarly for the other basic set constructions.

The reason why we prefer recursive subuniverses to inductive ones is that we obtain in inductive subuniverses many copies of the same element. In the simple situation, we have for instance two codes for N in V , namely \widehat{N}_V and $\widehat{T}_V(\widehat{N}_U)$. This is not very aesthetic (it doesn't cause many problems when actually making use of such universes). However, one should be aware that we anyway get often doubling of codes. For instance, if in case of the super universe to be introduced later in this article we form a recursive subuniverse (U, \widehat{T}_U) of V, T_V containing the family $\langle N_V, (x)\widehat{N}_V \rangle$, then U will have infinitely many codes for N , namely $\widehat{N}_U, \widehat{a}$ and $\widehat{b}(n)$ for $n : N$.

When modelling subuniverses, we will treat them as mere subsets of the interpretation of the original universe. \widehat{T}_U will for both the recursive and inductive subuniverses be essentially the identity function, expressed by the reduction rule $\widehat{T}_U(a) \longrightarrow a$. This means that in the model we will identify codes coming from the same element in U , so for instance all the codes for the natural numbers in the above example will be equal in the model.

3.2 Modelling Universes

For modelling universes closed under W in extensions of $KP\omega$ one needs the notion of a recursively inaccessible. A recursively inaccessible ordinal I is an admissible ordinal I s.t. for every $\beta \in I$ there exists a $\gamma \in I$ s.t. $\beta \in \gamma$ and γ is an admissible ordinal. If one refers to admissible sets rather than admissible ordinals, then a recursively inaccessible set is a set a s.t.

$$\text{Ad}(a) \wedge \forall y \in a. \exists z \in a. y \in z \wedge \text{Ad}(z) .$$

Let KPI^+ be the theory of $KP\omega$ extended by the existence of one recursively inaccessible set and finitely many admissibles above it. So KPI^+ has constants $a_{I,k}$ for Meta $k \in \omega$, axioms expressing that $a_{I,0}$ is a recursively inaccessible set and for Meta $k \in \omega$ the axioms

$$\text{Ad}(a_{I,k}) \wedge a_{I,k} \in a_{I,k+1} .$$

Proof theoretically equivalent is the theory KPI plus the existence of one recursively inaccessible set.

A universe is a family of sets, and therefore it will be modelled by a family of terms $\llbracket \text{Fam} \rrbracket(\text{set})$, as given by the following definition:

Definition 3.1. (in the language of Kripke-Platek set theory).

(a) By set we mean the class of all sets.

(b) $\llbracket \text{Fam} \rrbracket(\text{set}) := \{(A, B) \in \text{set}^2 \mid A \subseteq \llbracket \text{Term} \rrbracket^2$
 $\wedge B$ is a function
 $\wedge \text{dom}(B) = \text{Flat}(A)$
 $\wedge \forall x \in \text{Flat}(A). B(x) \subseteq \llbracket \text{Term} \rrbracket^2$
 $\wedge \forall \langle x, x' \rangle \in A. B(x) = B(x')\}$.

(c) If $(A, B) \in \llbracket \text{Fam} \rrbracket(\text{set})$. Then

$$\llbracket \text{Fam} \rrbracket(A, B) := (x \in A) \llbracket \times \rrbracket (B(x) \llbracket \rightarrow \rrbracket A) ,$$

(which is

$$\text{Clos}_{\rightarrow}(\{\langle c, c' \rangle \in \llbracket \text{Term} \rrbracket^2 \mid \langle \pi_0(c), \pi_0(c') \rangle \in A$$

$$\wedge (B(\pi_0(c)) \llbracket \rightarrow \rrbracket A) = (B(\pi_0(c')) \llbracket \rightarrow \rrbracket A)$$

$$\wedge \forall \langle x, x' \rangle \in B(\pi_0(c)). \langle \pi_1(c)(x), \pi_1(c')(x') \rangle \in A\}) .$$

(d) If $(A, B), (A', B') \in \llbracket \text{Fam} \rrbracket(\text{set})$. Then

$$(A, B) \leq (A', B') :\Leftrightarrow A \subseteq A' \wedge B \subseteq B' .$$

Note that by $(A, B), (A', B') \in \llbracket \text{Fam} \rrbracket(\text{set})$, $B \subseteq B'$ is equivalent to

$$B \upharpoonright \text{Flat}(A) = B' \upharpoonright \text{Flat}(A) .$$

(e) If $(U^\alpha, T^\alpha) \in \llbracket \text{Fam} \rrbracket(\text{set})$ for $\alpha < \beta$ s.t.

$$\forall \alpha < \alpha' < \beta. (U^\alpha, T^\alpha) \leq (U^{\alpha'}, T^{\alpha'}) .$$

Then we define

$$(U^{<\beta}, T^{<\beta}) := \left(\bigcup_{\alpha < \beta} U^\alpha, \bigcup_{\alpha < \beta} T^\alpha \right) .$$

Note that, if $x \in \text{Flat}(U^\alpha)$, $\alpha < \beta$, then

$$T^{<\beta}(x) = T^\alpha(x) .$$

A model of a universe (U, T) can be obtained by defining sets $\llbracket U \rrbracket$ and $\llbracket T \rrbracket$ s.t.

$$(\llbracket U \rrbracket, \llbracket T \rrbracket) \in \llbracket \text{Fam} \rrbracket(\text{set}) .$$

Then

$$\llbracket T(s) \rrbracket_\rho := \llbracket T \rrbracket(\llbracket s \rrbracket_\rho) .$$

We will introduce first by induction on α sets

$$(U^\alpha, T^\alpha) \in \llbracket \text{Fam} \rrbracket(\text{set})$$

s.t.

$$\alpha < \beta \rightarrow (U^\alpha, T^\alpha) < (U^\beta, T^\beta) .$$

then we define (assuming the existence of a recursively inaccessible ordinal I)

$$\llbracket U \rrbracket = U^{<I} , \quad \llbracket T \rrbracket = U^{<I}$$

for the least recursively inaccessible ordinal I . (In case we close the universe under additional constructions I will be replaced by a recursively inaccessible ordinal, which has some additional closure properties.)

We will guarantee that

$$\text{PER}(U^\alpha) \wedge \forall x \in \text{Flat}(U^\alpha) \text{PER}(T^\alpha(x)) .$$

We will as well enforce

$$U^\alpha \in L'_{\alpha+1} \wedge \forall x \in \text{Flat}(U^\alpha) T^\alpha(x) \in L'_\alpha .$$

All approximations of interpretations of universes U^α introduced in this article will be closed under reductions, which means in the following:

- If $\langle b, b' \rangle \in U^\alpha$, $c \longrightarrow^* b$, $c' \longrightarrow^* b'$ then $\langle c, c' \rangle \in U^\alpha$.
- If $b \longrightarrow^* c \in U^\alpha$, then $T^\alpha(b) = T^\alpha(c)$.

More formally the definition of U^α below has to be replaced by the definition of a set U'^α and function T'^α , and one defines then

$$\begin{aligned} U^\alpha &:= \text{Closure}(U'^\alpha) , \\ T^\alpha(a) &:= T'^\alpha(a') \text{ if } a \longrightarrow^* a' \in U'^\alpha . \end{aligned}$$

We will however in the following not explicitly mention U'^α , T'^α in this article, assuming that the reader can adapt the model correspondingly.

U^α is the closure of $U^{<\alpha}$ under one application of each basic set constructions of type theory, provided the results are in L'_α :

- $\llbracket N \rrbracket \in L'_\alpha \rightarrow \langle \widehat{N}, \widehat{N} \rangle \in U^\alpha$, $T^\alpha(\widehat{N}) = \llbracket N \rrbracket$.
- If $\langle \langle a, b \rangle, \langle a', b' \rangle \rangle \in \llbracket \text{Fam} \rrbracket(U^{<\alpha}, T^{<\alpha})$, and

$$A := \llbracket \Sigma \rrbracket(T^{<\alpha}(a), \lambda x \in T^{<\alpha}(a). T^{<\alpha}(b(x))) \in L'_\alpha ,$$

then $\langle \widehat{\Sigma}(a, b), \widehat{\Sigma}(a', b') \rangle \in U^\alpha$ and $T^\alpha(\widehat{\Sigma}(a, b)) = A$.

- Similarly for the other standard constructions of type theory.

In order to get a model of universe closed under some other principles, one will later add as well to U^α elements corresponding to those sets. In this article, these closure principles will always be the closure under certain universes.

We show that $(U^{<I}, T^{<I})$ is a model of the standard universe as follows (The proof shows in general that, if we close U^α under additional constructions in such a way that $U^\alpha \in L'_{\alpha+1}$ and for $x \in \text{Flat}(U^\alpha)$, $T^\alpha(x) \in L'_\alpha$, and if κ is a recursively inaccessible, then $(U^{<\kappa}, T^{<\kappa})$ is closed under the basic set constructions):

First one observes that if $A \subseteq \llbracket \text{Term} \rrbracket^2$, $A \in L'_\alpha$, for $a \in \text{Flat}(A)$, $B(x) \in \llbracket \text{Term} \rrbracket^2$ and $B(x) \in L'_\alpha$, then $\llbracket \Sigma \rrbracket(A, B)$, $\llbracket \Pi \rrbracket(A, B) \in L'_{\alpha+2}$ and $\llbracket W \rrbracket(A, B) \in L'_{\alpha+1}$. Similar properties hold for $\llbracket + \rrbracket$ and $\llbracket \text{Id} \rrbracket$ (making use of $\alpha + 2$).

Now, if I is a recursively inaccessible ordinal, then $(U^{<I}, T^{<I})$ is closed under $\widehat{N}_n, \widehat{N}, \Sigma, \Pi, \widehat{W}, \widehat{+}, \widehat{\text{Id}}$: We consider only the most difficult case \widehat{W} : First we show a small Lemma, which will be useful in general and will hold for all universe constructions defined in this article:

Lemma 3.2. *Let κ be admissible, and*

$$\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket(U^{<\kappa}, T^{<\kappa}) .$$

then there exists an $\alpha < \kappa$ s.t.

$$\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket(U^{<\alpha}, T^{<\alpha}) .$$

Proof of the Lemma: We have $\langle a, a' \rangle \in U^{<\kappa}$, therefore there exists $\alpha_0 < \kappa$ s.t. $\langle a, a' \rangle \in U^{<\alpha_0}$. Furthermore, $T^{<\alpha_0}(a) \in L'_{\alpha_0}$, and for all $\langle x, x' \rangle \in T^{<\alpha_0}(a)$ there exists $\alpha_1 < \kappa$ s.t. $\langle b(x), b'(x') \rangle \in U^{<\alpha_1}$. By κ admissible we can find a uniform α_1 for this. Now let $\alpha := \max\{\alpha_0, \alpha_1\}$.

We continue with the main proof. Assume $\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket(U^{<I}, T^{<I})$. Then by Lemma 3.2 we find an α s.t. $\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket(U^{<\alpha}, T^{<\alpha})$. We have $T^{<I}(a) = T^{<\beta}(a) \in L'_\beta$, $T^{<I}(b(x)) = T^{<\beta}(b(x)) \in L'_\beta$, therefore $\llbracket W \rrbracket(T^{<I}(a), \lambda x \in T^{<I}(a). T^{<I}(b(x))) \in L'_{\beta+1}$. Therefore $\langle\widehat{W}(a, b), \widehat{W}(a', b')\rangle \in U^{\beta+2} \subseteq U^{<I}$.

One obtains a model of the type theory ML_WU in KPI^+ as follows: We have modelled all the ingredients already, and need only to make sure that the constructions can be carried out in KPI^+ . Let κ_n be the n th admissible above I . First, one observes, that $\llbracket U \rrbracket \in L'_{I+2}$, for $x \in \text{Flat}(\llbracket U \rrbracket)$ we have $\llbracket T_U \rrbracket(x) \in L'_I$, and for $\langle x, x' \rangle \in \llbracket U \rrbracket$ we have $\llbracket T_U \rrbracket(x) = \llbracket T_U \rrbracket(x')$. Therefore one defines $o(U) := I + 2$, $o(T_U(s)) := I$. For all other set constructions we have seen that if the maximum of the ordinals of the subterms is α , then the ordinal of the term itself is $\alpha + 2$ or $\alpha^+ + 1$. Therefore we obtain that for any set term A there exists an n s.t. $o(A) \leq \kappa_n$ and therefore for all ρ $\llbracket A \rrbracket_\rho \in L'_{\kappa_n}$. One sees now that we can interpret all sets in KPI^+ , and can show that this model is sound with respect to the above mentioned correctness conditions. So we have given a sketch of the following theorem of which a detailed proof can be found in [25]; see as well [27]; Griffor and Rathjen have shown related results in [14]:

Theorem 3.3.

- (a) ML_{WU} can be modelled in KPI^+ .
- (b) $|\text{ML}_{\text{WU}}| \leq |\text{KPI}^+|$.
- (c) The previous statements hold as well if we replace intensional by extensional equality.

In the other direction, we have developed in [29] (see as well our PhD thesis [25], the overview articles [28] and [31]; Rathjen and Griffor have obtained related results in [14]) a well-ordering proof which shows the other direction, namely:

Theorem 3.4. (a) $|\text{ML}_{\text{WU}}| = |\text{KPI}^+|$

- (b) The previous statements hold as well if we replace intensional by extensional equality.

4 Palmgren's Super Universe and Recursively Hyperinaccessibles

4.1 Definition of the Super Universe

The first substantial step beyond the type theory ML_{WU} was the introduction of a hierarchy of finitely many universes by Per Martin-Löf. We won't cover this extension in this article, and will instead move to the next step, namely the super universe, as introduced by Erik Palmgren. He introduced a super universe operator, which forms a universe containing a family of sets. The super universe is defined as a universe closed under this operator.

In order to introduce this in detail, let us first introduce some notions:

Definition 4.1. (Refers to the logical framework).

- (a) The type of families of sets is defined (using the logical framework) as $\text{Fam}(\text{Set}) := (X : \text{Set}) \times (X \rightarrow \text{Set})$.
- (b) Let $\langle U, T \rangle \in \text{Fam}(\text{Set})$. The family of sets in $\langle U, T \rangle$ is the set of families of sets in U , indexed over elements of U : $\text{Fam}(U, T) := (x : U) \times (T(x) \rightarrow U)$.
- (c) Let $\langle U, T \rangle \in \text{Fam}(\text{Set})$. The lifting of an element $a : \text{Fam}(U, T)$ to an element of $\text{Fam}(\text{Set})$ is defined as $T[a] := \langle T(\pi_0(a)), (x)T(\pi_1(a)(x)) \rangle : \text{Fam}(\text{Set})$.

Palmgren, first introduced the super universe operator

$$\text{SU} : \text{Fam}(\text{Set}) \rightarrow \text{Fam}(\text{Set}) .$$

$\text{SU}(A)$ is a universe $\langle U', T' \rangle$ closed under the standard set constructions and containing A . This means that if $A = \langle A', B' \rangle$, then U' contains

- a code \hat{a} for A' , i.e. $T'(\hat{a}) = A'$,
- and for $a : A'$ contains a code $\hat{b}(a)$ for $B'(a)$, i.e. $T'(\hat{b}(a)) = B'(a)$.

The next step was to introduce the super universe as a universe (V, T_V) , which is as usual closed under the standard set constructions and which reflects the super universe operator: if $a : \text{Fam}(V, T_V)$, then V contains codes $\langle \hat{U}_a, (x)\hat{T}_a(x) \rangle$ for $\text{SU}(T_V[a])$, which means that

$$T_V[\langle \hat{U}_a, (x)\hat{T}_a(x) \rangle] = \text{SU}(T_V[a]) .$$

We will in the following present a version, in which $(T(\hat{U}_a), \hat{T}_a)$ is a recursive subuniverse of (V, T_V) . Furthermore, we will define the universe operator only restricted to elements of $\text{Fam}(V, T_V)$ (which will be written as $U_{a,b}$ for $\langle a, b \rangle : \text{Fam}(V, T_V)$). We uncurry SU and obtain the following rules:

$$\begin{array}{c} V : \text{Set} \\ \hline a : V \quad b : T_V(a) \rightarrow V \\ \hline U_{a,b} : \text{Set} \\ \hline a : V \quad b : T_V(a) \rightarrow V \\ \hline \hat{U}_{a,b} : V \\ \hline a : V \quad b : T_V(a) \rightarrow V \quad c : U_{ab} \\ \hline \hat{T}_{U,a,b}(c) : V \end{array} \quad \begin{array}{c} \frac{a : V}{T_V(a) : \text{Set}} \\ \\ \\ \\ \\ \\ T_V(\hat{U}_{a,b}) = U_{a,b} \end{array}$$

Let $T_{U,a,b}(x) := T_V(\hat{T}_{U,a,b}(x)) : \text{Set}$.

(V, T_V) is a universe, which means it is closed under the basic set constructions. The codes are written as $\hat{N}_V, \hat{\Sigma}_V$, etc.

One demands for $a : V, b : T_V(a) \rightarrow V$ that $(U_{a,b}, \hat{T}_{U,a,b})$ is a recursive subuniverse of (V, T_V) , which means as indicated in Sect. 3 that $U_{a,b}$ is closed under constructors for forming codes for the standard set constructions like $\hat{N}_{U,a,b}, \hat{\Sigma}_{U,a,b}$, etc., and that $\hat{T}_{U,a,b}$ decodes them as the corresponding codes in V , e.g. $\hat{T}_{U,a,b}(\hat{N}_{U,a,b}) = \hat{N}_V$.

Furthermore, $U_{a,b}$ contains codes for a, b , which means that one has the following rules:

$$\begin{array}{c} \hat{a}_{a,b} : U_{a,b} \quad \hat{T}_{U,a,b}(\hat{a}_{a,b}) = a \\ \\ \frac{c : T_{U,a,b}(a)}{\hat{b}_{a,b}(c) : U_{a,b}} \quad \hat{T}_{U,a,b}(\hat{b}_{a,b}(c)) = b(c) \end{array}$$

Let $\text{ML}_W + (\text{Superuniv})$ be the type theory consisting of the rules of the small logical framework, the basic set constructions, and the rules for $V, T_V, U_{ab}, \hat{U}_{ab}, \hat{T}_{U,a,b}$ as above.

In this type theory, one can define using the elimination rules for \mathbb{N} a hierarchy of universes $\widehat{U}'_n : \mathbb{V}$ by

$$\widehat{U}'_0 = \widehat{U}_{\widehat{N}_0, (x)\widehat{N}_0} \quad \widehat{U}'_{\text{TV}(n)} = \widehat{U}_{\widehat{N}_1, (n)\widehat{U}'_n}$$

but as well a universe containing this hierarchy of universes, namely $U_{\widehat{N}, (n)\widehat{U}'_n}$.

This is the first step towards forming a transfinite hierarchy of universes: for any well-founded relation (with well-foundedness provable in this type theory) R, \prec we can define a hierarchy of universes \widehat{U}''_α for $\alpha \in R$ s.t. \widehat{U}''_α contains codes for \widehat{U}''_β for $\beta \prec \alpha$.

4.2 A Model of the Super Universe

$\text{ML}_W + (\text{Superuniv})$ has strength $(\text{KP} + (\text{hyper} - \text{inacc}))^+$, where $(\text{KP} + (\text{hyper} - \text{inacc}))^+$ is Kripke-Platek theory plus the existence of one recursively hyper-inaccessible set plus finitely many admissibles above it, which is equivalent to KPI plus the existence of one recursively hyper-inaccessible set.

Here a recursively hyper-inaccessible set is an admissible set a s.t. for all $x \in a$ there exists a $y \in a$ s.t. $x \in y$ and y is recursively inaccessible. A recursively hyper-inaccessible ordinal is an admissible ordinal HI s.t.

$$\forall \alpha \in \text{HI}. \exists \beta \in \text{HI}. \alpha \in \beta \wedge \text{Inacc}(\beta) ,$$

where $\text{Inacc}(\beta)$ means that β is a recursively inaccessible ordinal.

The precise formulation of $(\text{KP} + (\text{hyper} - \text{inacc}))^+$ is similar to KPI^+ , except that one replaces “inaccessible” by “hyper-inaccessible”.

An upper bound for the strength of $\text{ML}_W + (\text{Superuniv})$ can be obtained by modelling it in $(\text{KP} + (\text{hyper} - \text{inacc}))^+$. We will give in the following a sketch of the model, full details will be presented in a future paper.

The basic construction is outlined in Sect. 2. As pointed out there, in $\llbracket \text{Term} \rrbracket$ we forget about the constants used for type checking only, so we have for instance only a constant \widehat{N} and interpret

$$\llbracket \widehat{N}_V \rrbracket_\rho = \llbracket \widehat{N}_{U, a, b} \rrbracket_\rho = \widehat{N} ,$$

similarly for $\widehat{\Sigma}, \widehat{a}, \widehat{b}$. $\widehat{U}, \widehat{T}_U$ will depend on a, b , and we will add the reduction rules for \widehat{T}_V , expressed by the above, i.e.

$$\begin{aligned} \widehat{T}_{U, a, b}(\widehat{N}) &\longrightarrow \widehat{N} , \\ \widehat{T}_{U, a, b}(\widehat{\Sigma}(c, d)) &\longrightarrow \widehat{\Sigma}(c, d) , \\ \widehat{T}_{U, a, b}(\widehat{a}) &\longrightarrow a , \\ \widehat{T}_{U, a, b}(\widehat{b}(c)) &\longrightarrow b(c) . \end{aligned}$$

As for standard universes, the super universe is interpreted, by defining first by induction over α sets

$$(\mathbb{V}^\alpha, \mathbb{T}^\alpha) \in \llbracket \text{Fam} \rrbracket(\text{set})$$

s.t.

$$\alpha < \beta \rightarrow (U_\alpha, T_\alpha) < (U_\beta, T_\beta) .$$

As for all models of universes in this article, we add to V^α the codes for standard set constructors applied to arguments, if the elements of the universe they refer to are already in $V^{<\alpha}$ and the sets added are in L'_α . Additionally, whenever there exists a β s.t. $\beta+1 < \alpha$ and $V^{<\beta}$ is closed under all universe constructions, and s.t.

$$\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket (V^{<\beta}, T^{<\beta}) ,$$

then

$$\langle \widehat{U}_{a,b}, \widehat{U}_{a',b'} \rangle \in V^\alpha$$

and we define

$$T^\alpha(\widehat{U}_{a,b}) := V^{<\beta}$$

for the minimal β , which has the mentioned closure properties.

Note that with this construction we obtain that

$$\begin{aligned} V^\alpha &\in L'_{\alpha+1} , \\ \forall x \in \text{Flat}(V^\alpha). T^\alpha(x) &\in L'_\alpha . \end{aligned}$$

Let HI be the least recursively hyper-inaccessible. Then one interprets

$$\begin{aligned} \llbracket V \rrbracket &:= V^{<\text{HI}} , \\ \llbracket T_V \rrbracket &:= T^{<\text{HI}} , \\ \llbracket U_{a,b} \rrbracket_\rho &:= T^{<\text{HI}}(\llbracket \widehat{U}_{a,b} \rrbracket_\rho) . \end{aligned}$$

Note that this model construction won't admit elimination rules for $U_{a,b}$, since it is not interpreted as the least set closed under a, b and the basic set constructions – the $V^{\beta'}$ chosen above might contain elements of the form $\widehat{U}_{a',b'}$, which have been added, since the model of V needs to be closed under \widehat{U} , but which are not elements of the least subuniverse closed under a, b and the standard set constructions. One could modify the above model in order to admit elimination rules: Then one would try to form at each stage α the least subset of $V^{<\alpha}$ closed under the universe constructions and containing a and $b(x)$ for $x \in \text{Flat}(T^{<\alpha}(a))$. If this definition succeeds and the least such set is A , one would add $\widehat{U}_{a,b}$ to V^α and define $T^\alpha(\widehat{U}_{a,b}) := A$.

We have to verify that $V^{<\text{HI}}$ has the necessary closure properties of the super universes. Since HI is a recursively inaccessible, $V^{<\text{HI}}$ will be closed under universe constructions as before. $V^{<\text{HI}}$ is as well closed under $(a, b)\widehat{U}_{a,b}$: Assume

$$\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket (V^{<\text{HI}}, T^{<\text{HI}}) .$$

By the admissibility of HI it follows by Lemma 3.2

$$\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket (V^{<\alpha}, T^{<\alpha})$$

for some $\alpha < \text{HI}$ Since HI is a recursively hyper inaccessible, there exists a recursively inaccessible ordinal $\alpha < \kappa < \text{HI}$. V^κ will be closed under the universe constructions and

$$\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket (V^\kappa, T^\kappa) ,$$

therefore

$$\langle \widehat{U}_{a,b}, \widehat{U}_{a',b'} \rangle \in V^{\kappa+2} \subseteq V^{< \text{HI}} .$$

That $\llbracket U_{a,b} \rrbracket_\rho$ has the closure properties needed in order to interpret $U_{a,b}$ follows by the construction.

The remaining construction is as for any standard model of type theory, i.e. we interpret the constructions of the small logical framework and the basic set constructions on top of U . As for $\text{ML}_W U$, we need $< \omega$ many admissibles in order to interpret W -types built on top of $V^{< \text{HI}}$, therefore finite approximations of the type theory can be interpret into Kripke Platek set theory plus one recursively hyper-inaccessible and finitely many admissibles above it. So we have given the essence of a proof of the following theorem:

Theorem 4.2.

- (a) *We can model $\text{ML}_W + (\text{Superuniv})$ in $(\text{KP} + (\text{hyper} - \text{inacc}))^+$.*
- (b) $|\text{ML}_W + (\text{Superuniv})| \leq |(\text{KP} + (\text{hyper} - \text{inacc}))^+|$
- (c) *The previous statements hold as well if we replace intensional by extensional equality.*

In order to obtain the precise proof theoretic strength, one has to carry out a well-ordering proof in order to obtain

$$|\text{ML}_W + (\text{Superuniv})| \geq |(\text{KP} + (\text{hyper} - \text{inacc}))^+| .$$

We note that if one adds as in Erik's original type theory the full universe operator $\text{SU} : \text{Fam}(\text{Set}) \rightarrow \text{Fam}(\text{Set})$, we get a little bit more strength. Formulated without the logical framework, in this type theory one has depending on $A : \text{Set}$ and $x : A \Rightarrow B : \text{Set}$ the set $U(A, (x)B) : \text{Set}$ and decoding function

$$u : U(A, (x)B) \Rightarrow T(A, (x)B, u) : \text{Set} .$$

$(U(A, (x)B), T(A, (x)B))$ is a universe containing codes for A and $B[x := a]$, and in the above definition one replaces

$$\begin{array}{ll} U_{a,b} & \text{by } U(T_V(a), (x)T_V(b(x))) , \\ T_{a,b}(c) & \text{by } T(T_V(a), (x)T_V(b(x))) . \end{array}$$

In this theory one can apply the super universe operator as well to V itself, and to elements formed from it. However, only finite nestings of the super universe operator are possible. The modelling of one application of the super universe to arguments (A, B) requires one recursively inaccessible above $\max\{o(A), o(B)\}$. The modelling of V requires one recursively hyper-inaccessible, so in total one needs for Palmgren's original type theory one recursively hyper-inaccessible ordinal and finitely many recursively inaccessible ordinals above it.

5 The Mahlo Universe and Extensions of It

5.1 Definition of the Mahlo Universe

Once one has defined the super universe, one can move to a hyper-universe by forming a universe U^2 with decoding function T^2 , s.t. for any family of sets a, b in U^2 there exists a super universe $U_{a,b}^1$ and an element $\widehat{U}_{a,b}^1$ in V , s.t. $T_V(\widehat{U}_{a,b}^1) = U_{a,b}^1$, where $U_{a,b}^1$ contains codes which decode as a, b as before. So we have a decoding function $\widehat{T}_{a,b}^1 : U_{a,b}^1 \rightarrow U^2$, and define $T_{a,b}^1 := T^2 \circ \widehat{T}_{a,b}^1$. $U_{a,b}^1$ is closed under the standard universe constructions. That $U_{a,b}^1$ is a super universe means, that for any family of sets c, d in $U_{a,b}^1$ there exists a subuniverse $U_{a,b,c,d}^2$ of $U_{a,b}^1$ which contains codes decoding as c, d , is closed under the standard universe constructions, and which is represented as a code in $U_{a,b}^1$.

One could define now hyper $^\alpha$ -super universes, but we won't spell this out in this article. The construction is similar to the step from a Mahlo-universe to a hyper $^\alpha$ -Mahlo universe, which will be discussed in the followup of this article [32].

We note that E. Palmgren has [20] introduced higher order universes. In those universes, the lowest level is a universe, the next level is a universe of universe operators, the next level is a universes of operators on universe operators, etc. E. Palmgren added rules expressing the closure of higher order universes under the application of a universe operator from the next level. It is conjectured that E. Palmgren's universe reaches the strength of KPM, which is slightly below the strength of the Mahlo universe. We will not investigate this construction any further (see the remark in the introduction that Palmgren's universe contains a Mahlo universe in disguise, even so it is not easy to see this).

We will investigate in which sense the hyper $^\alpha$ -super universes are special instances of universes closed under operators on families of sets, and in which sense the Mahlo universe to be introduced generalises this.

Definition 5.1. (*Assumes the logical framework*).

- (a) An operator on families of sets is a function $f : \text{Fam}(\text{Set}) \rightarrow \text{Fam}(\text{Set})$.
- (b) Let $\langle U, T \rangle : \text{Fam}(\text{Set})$. An operator on families of elements of U is a function $f : \text{Fam}(U, T) \rightarrow \text{Fam}(U, T)$.
- (c) Let $\langle A, B \rangle, \langle A', B' \rangle : \text{Fam}(\text{Set})$. Assume elimination of $+$ into Set . Then

$$\langle A, B \rangle \cup_{\text{Fam}(\text{Set})} \langle A', B' \rangle := \langle C, D \rangle$$

where

$$C := (A + N_1) + (A' + N_1)$$

and D is such that (remember that A_0^1 is the canonical element of N_1)

$$\begin{aligned} D(\text{inl}(\text{inl}(a))) &= B(a) \quad , \quad D(\text{inl}(\text{inr}(A_0^1))) = A \quad , \\ D(\text{inr}(\text{inl}(a'))) &= B'(a') \quad , \quad D(\text{inr}(\text{inr}(A_0^1))) = A' \quad . \end{aligned}$$

The super universe operator SU takes $A : \text{Fam}(\text{Set})$ and forms a universe closed under the operator $(X)A : \text{Fam}(\text{Set}) \rightarrow \text{Fam}(\text{Set})$. The super universe is a universe which is closed under the super universe operator SU . One can define now a hyper-universe operator HU , which takes $A : \text{Fam}(\text{Set})$ and forms a super universe closed $(X)A$. One obtains essentially the same set by taking a universe closed under $(X)SU(X) \cup_{\text{Fam}(\text{Set})} A$ (assuming this is definable). Then one can form a hyper²-universe operator, a universe closed under it etc.

The step towards the Mahlo universe is now to form a universe, which is closed under arbitrary operators on families of elements of itself. This means that we have a universe (V, T_V) , s.t. for any operator $f : \text{Fam}(V, T_V) \rightarrow \text{Fam}(V, T_V)$ there exists a subuniverse (U_f, T_f) of V , closed under f and represented as an element \widehat{U}_f in V .

Spelled out without using the logical framework, we obtain the following rules:

- We have the rules of the small logical framework.
- We have the rules of the basic set constructions.
- We have rules expressing that (V, T_V) is a universe closed under the basic set constructions. The constructors for the codes for these set constructions will be denoted by $\widehat{N}_V, \widehat{\Sigma}_V$, etc.
- We have that for every operator f on families of sets in (V, T_V) there is a subuniverse U_f closed under it. We split f into two components and uncurry it, and obtain the following rules:

$$\frac{f : (a : V, b : T_V(a) \rightarrow V) \rightarrow V \quad g : (a : V, b : T_V(a) \rightarrow V, T_V(f(a, b))) \rightarrow V}{U_{f,g} : \text{Set}}$$

- We have rules expressing that $U_{f,g}$ is a recursive subuniverse closed under the standard set constructions. So we have the decoding function

$$\widehat{T}_{U,f,g} : U_{f,g} \rightarrow V$$

and define

$$T_{U,f,g}(a) := T_V(\widehat{T}_{U,f,g}(a)) : \text{Set} .$$

The codes for the basic set constructions will be denoted by $\widehat{N}_{U,f,g}, \widehat{\Sigma}_{U,f,g}$, etc.

- We have rules expressing that $(U_{f,g}, T_{f,g})$ is closed under f and g :

$$\begin{aligned} \widehat{f}_{f,g} & : (a : U_{f,g}, b : T_{U,f,g}(a) \rightarrow U_{f,g}) \rightarrow U_{f,g} \\ \widehat{T}_{U,f,g}(\widehat{f}_{f,g}(a, b)) & = f(\widehat{T}_V(a), \widehat{T}_V \circ b). \\ \widehat{g}_{f,g} & : (a : U_{f,g}, b : T_{U,f,g}(a) \rightarrow U_{f,g}, T_V(f(\widehat{T}_V(a), \widehat{T}_V \circ b))) \rightarrow U_{f,g} \\ \widehat{T}_{U,f,g}(\widehat{g}_{f,g}(a, b, c)) & = g(\widehat{T}_V(a), \widehat{T}_V \circ b, c). \end{aligned}$$

- And finally, we have the rule that V contains a code for $U_{f,g}$. As P. Martin-Löf observed this is the rule which makes the Mahlo universe so strong – having just subuniverses closed under f, g doesn't add any strength to $ML_W U$.

$$\frac{f : (a : V, b : T_V(a) \rightarrow V) \rightarrow V \quad g : (a : V, b : T_V(a) \rightarrow V, T_V(f(a, b))) \rightarrow V}{\widehat{U}_{f,g} : V} \\ \widehat{T}_V(\widehat{U}_{f,g}) = U_{f,g}$$

We call the resulting type theory $ML_W + (\text{Mahlo})$.

5.2 A Model of the Mahlo Universe

$ML_W + (\text{Mahlo})$ has strength $(KPM)^+$, where $(KPM)^+$ is $KP\omega$ plus the existence of one recursively Mahlo set plus finitely many admissibles above it, which is equivalent to KPI plus the existence of one recursively Mahlo set.

Here a *recursively Mahlo set* is an admissible set a_M s.t.

$$(\text{Mahlo}) \quad \forall \vec{z} \in a_M. ((\forall x \in a_M. \exists y \in a_M. \varphi(x, y, \vec{z})) \\ \rightarrow \exists z \in a_M. \text{Ad}(z) \wedge \vec{z} \in a_M. \wedge (\forall x \in z. \exists y \in z. \varphi(x, y, \vec{z}))) \\ \text{where } \varphi(x, y, \vec{z}) \text{ is } \Delta_0$$

A *recursively Mahlo ordinal* is an ordinal M s.t. L_M is a recursively Mahlo ordinal.

The basic construction of the Mahlo universe is as for the super universe: We define by recursion over α approximations V^α of $\llbracket V \rrbracket$ by closing it systematically under the basic set constructions, and by adding suitable codes $\widehat{U}_{f,g}$ to it.

The question is when to add codes $\widehat{U}_{f,g}$ to the Mahlo universe. In the introduction rules, f, g need to be total functions from families of elements in V into families of elements in V . But we only know this type of total functions, once the definition of $\llbracket V \rrbracket$ is complete. The trick to get around this problem is to add more elements $\widehat{U}_{f,g}$ to the universe than are actually justified by the introduction rules. We only demand that f, g are total on the subuniverse $U_{f,g}$ itself. So we take f, g to be arbitrary terms. Then we try to form a subuniverse of the current approximation of V^α of $\llbracket V \rrbracket$, which is closed under f and g . If we succeed, we add $\widehat{U}_{f,g}$ to $\llbracket V \rrbracket$ and define $\llbracket T_V \rrbracket(\widehat{U}_{f,g})$ as the least such subuniverse.

If we assume now total functions f, g from families of elements of $V^{<M}$ into families of elements of $V^{<M}$, then $V^{<M}$ is closed under f, g , where M is a recursively Mahlo ordinal. So there is at least one subuniverse of $V^{<M}$ closed under f, g , namely $V^{<M}$ itself. Then we will use the fact that M is a Mahlo-ordinal, and find that such a universe occurred already at some stage $\alpha < M$, and that therefore $\widehat{U}_{f,g}$ is in $V^{<M}$.

As for the super universe, we simplify this construction slightly and don't search for the least subuniverse closed under f and g , but only for the least β s.t. V^β is closed under f and g .

More precisely we proceed as follows: As in the model for the super universe, we forget in $\widehat{\mathbb{N}}$ the terms of the model the typing information needed in $\widehat{\mathbb{N}}$, $\widehat{\Sigma}$, etc. (but not in \widehat{f} , \widehat{g}), so for instance

$$\llbracket \widehat{\mathbb{N}}_V \rrbracket_\rho := \llbracket \widehat{\mathbb{N}}_{U,f,g} \rrbracket_\rho := \widehat{\mathbb{N}} .$$

We obtain essentially the same reduction rules for $\widehat{\mathbb{T}}_{U,f,g}$ as the reduction rules for $\widehat{\mathbb{T}}_{U,a,b}$ in case of the super universe. Furthermore we have the reductions

$$\begin{aligned} \widehat{f}_{f,g}(a,b) &\longrightarrow f(a,b) , \\ \widehat{g}_{f,g}(a,b,c) &\longrightarrow g(a,b,c) . \end{aligned}$$

As for the super universe we form $(V^\alpha, T^\alpha) \in \llbracket \text{Fam} \rrbracket(\text{set})$ by closing it under the standard set constructions. We will then interpret $\llbracket V \rrbracket = V^{<M}$, $\llbracket T_V \rrbracket = T^{<M}$.

Instead of closing it under $\widehat{U}_{a,b}$ for $\langle a, b \rangle \in \llbracket \text{Fam} \rrbracket(V^{<\alpha}, T^{<\alpha})$ we close it under $\widehat{U}_{f,g}$, where f, g form the two components of a function

$$\llbracket \text{Fam} \rrbracket(V^{<\alpha}, T^{<\alpha}) \llbracket \rightarrow \rrbracket \llbracket \text{Fam} \rrbracket(V^{<\alpha}, T^{<\alpha}) ,$$

and we call $\langle f, g \rangle$ an operator on families of sets in $(V^{<\alpha}, T^{<\alpha})$. More precisely, we introduce the following definition, where we first define a notion for the two components of an operator on families of sets in $\langle A, B \rangle$, and then define what it means to be an operator on families of such sets as the conjunction of the two:

Definition 5.2. *Let $\langle A, B \rangle \in \llbracket \text{Fam} \rrbracket(\text{set})$.*

$$(a) \llbracket \text{FamOper} \rrbracket_0(A, B) := (x \in A) \llbracket \rightarrow \rrbracket ((y \in B(x)) \llbracket \rightarrow \rrbracket A) \llbracket \rightarrow \rrbracket A .$$

(b) *Assume $f \in \text{Flat}(\llbracket \text{FamOper} \rrbracket_0(A, B))$. Then*

$$\llbracket \text{FamOper} \rrbracket_1(A, B, f) := (x \in A) \llbracket \rightarrow \rrbracket ((y \in B(x)) \llbracket \rightarrow \rrbracket A) \llbracket \rightarrow \rrbracket B(f(x, y)) \llbracket \rightarrow \rrbracket A .$$

(c) *The set of operators on families of sets in $\langle A, B \rangle$ is given as*

$$\llbracket \text{FamOper} \rrbracket(A, B) := \{ \langle \langle f, g \rangle, \langle f', g' \rangle \rangle \in (\llbracket \text{Term} \rrbracket^2)^2 \mid \langle f, f' \rangle \in \llbracket \text{FamOper} \rrbracket_0(A, B) \wedge \langle g, g' \rangle \in \llbracket \text{FamOper} \rrbracket_1(A, B, f) \} .$$

With this notation, the addition of $\widehat{U}_{f,g}$ to V^α is defined as follows:

Assume $f, f', g, g' \in \llbracket \text{Term} \rrbracket$. Assume $\alpha < M$ and β s.t. $\beta + 1 < \alpha$ and such that $V^{<\beta}$ is closed under all universe constructions. Assume that

$$\langle \langle f, g \rangle, \langle f', g' \rangle \rangle \in \llbracket \text{FamOper} \rrbracket(V^{<\beta}, T^{<\beta})$$

Let β be minimal that this property holds. Then

$$\langle \widehat{U}_{f,g}, \widehat{U}_{f',g'} \rangle \in V^\alpha ,$$

and define

$$T^\alpha(\widehat{U}_{f,g}) := V^{<\beta} .$$

Note that this construction preserves the property that

$$V^\alpha \in L'_{\alpha+1} \wedge \forall a \in \text{Flat}(V^\alpha). T^\alpha(a) \in L'_\alpha .$$

Now define

$$\begin{aligned} \llbracket U \rrbracket_\rho &:= V^{<M} , \\ \llbracket T_V(c) \rrbracket_\rho &:= T^{<M}(\llbracket c \rrbracket_\rho) , \\ \llbracket U_{f,g} \rrbracket_\rho &:= T^{<M}(\widehat{U}_{\llbracket f \rrbracket_\rho, \llbracket g \rrbracket_\rho}) . \end{aligned}$$

where M is the least recursively Mahlo ordinal.

We note that $T^\alpha(\widehat{U}_{f,g})$ is not defined as a minimal set closed under universe constructions, and under f and g : We could have done so by modifying the above definition and trying at stage α to define the least subset of V^α closed under f and g and the standard universe constructions. If this succeeds, this subset is independent of α (provided α is big enough so that this subset lies inside V^α) and we could have interpreted $T^\alpha(\widehat{U}_{f,g})$ as this set. This would interpret elimination rules for $U_{f,g}$ as well. But since elimination rules for $U_{f,g}$ don't add any proof theoretic strength, and that model construction is more complicated, we use the original approach.

We have to verify that $V^{<M}$ has the necessary closure properties of the Mahlo universe. Every recursively Mahlo ordinal is recursively inaccessible, therefore $V^{<M}$ will be closed under universe constructions as before. $V^{<M}$ is as well closed under $(f, g)\widehat{U}_{f,g}$: Assume $f, f', g, g' \in \llbracket \text{Term} \rrbracket$. Assume

$$\begin{aligned} \langle f, f' \rangle &\in \llbracket (x : V, y : T_V(x) \rightarrow V) \rightarrow V \rrbracket , \\ \langle g, g' \rangle &\in \llbracket (x : V, y : T_V(x) \rightarrow V, z(T_V(x), T_V \circ y)) \rightarrow V \rrbracket_{[z=f]} . \end{aligned}$$

This means

$$\langle \langle f, g \rangle, \langle f', g' \rangle \rangle \in \llbracket \text{FamOper} \rrbracket(V^{<M}, T^{<M})$$

We need to find a $\kappa < M$ s.t. $V^{<\kappa}$ is closed under universe constructions and that the previous conditions hold with M replaced by κ . For being closed under universe constructions it suffices that κ is inaccessible. We show that for every $\beta < M$ there exists a $\rho < M$ s.t. ρ is admissible, $\beta < \rho$ and s.t.

$$\begin{aligned} \langle \langle a, b \rangle, \langle a', b' \rangle \rangle &\in \llbracket \text{Fam} \rrbracket(V^{<\beta}, T^{<\beta}) \\ \rightarrow \langle \langle f(a, b), g(a, b) \rangle, \langle f'(a', b'), g'(a', b') \rangle \rangle &\in \llbracket \text{Fam} \rrbracket(V^{<\rho}, T^{<\rho}) \end{aligned}$$

This formula will then be expressed by a Π_2 -formula and by the Mahloness of M we then find a κ s.t. for $\beta < \kappa$ there exists a $\rho < \kappa$ which fulfils the previous conditions. Then κ will be inaccessible and $V^{<\kappa}$ will be closed under $\langle f, f' \rangle$ and $\langle g, g' \rangle$.

So assume $\beta < M$ and

$$\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket(V^{<\beta}, T^{<\beta}) .$$

Then there exists a γ s.t. $\langle f(a, b), f'(a', b') \rangle \in V^{<\gamma}$. Furthermore,

$$T^{<M}(f(a, b)) = T^{<\gamma}(f(a, b)) \in L'_\gamma ,$$

and for $\langle c, c' \rangle \in T^{<M}(f(a, b))$ there exists a $\delta < M$ s.t.

$$\langle g(a, b, c), g'(a', b', c') \rangle \in T^{<\delta}(f(a, b)) .$$

By the admissibility of M it follows that there exists a $\rho' < M$ s.t. for $\langle c, c' \rangle \in T^{<M}(f(a, b))$ we have

$$\langle g(a, b, c), g'(a', b', c') \rangle \in T^{<\rho'}(f(a, b)) .$$

We can obtain that $\gamma \leq \rho'$ and therefore

$$\langle\langle f(a, b), f'(a', b') \rangle, \langle g(a, b), g'(a', b') \rangle\rangle \in \llbracket \text{Fam} \rrbracket(V^{<\rho'}, T^{<\rho'}) .$$

Using again the admissibility of M and the fact that $\llbracket \text{Fam} \rrbracket(V^{<\rho'}, T^{<\rho'}) \in L'_{\rho'+2}$ we obtain depending on β a uniform ρ s.t. if

$$\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket(V^{<\beta}, T^{<\beta})$$

then

$$\langle\langle f(a, b), f'(a', b') \rangle, \langle g(a, b), g'(a', b') \rangle\rangle \in \llbracket \text{Fam} \rrbracket(V^{<\rho}, T^{<\rho}) .$$

Furthermore, using the fact that every recursively Mahlo ordinal is recursively inaccessible, we can achieve $\beta < \rho$ admissible. Therefore, for every $\beta < M$ there exists a $\beta < \rho < M$ s.t. ρ is admissible and s.t. if

$$\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket(V^{<\beta}, T^{<\beta}) ,$$

then

$$\langle\langle f(a, b), g(a, b) \rangle, \langle f'(a', b'), g'(a', b') \rangle\rangle \in \llbracket \text{Fam} \rrbracket(V^{<\rho}, T^{<\rho}) .$$

The property we have shown can be expressed as a Π_2 -formula, i.e. a formula of the form

$$\forall \beta < M. \exists \rho < M. \varphi(\beta, \rho) ,$$

where φ is Δ_0 .

By the Mahlo property we obtain the existence of an admissible $\kappa < M$ s.t.

$$\forall \beta < \kappa. \exists \rho < \kappa. \varphi(\beta, \rho) .$$

Since ρ is always an admissible $> \beta$, it follows that κ is recursively inaccessible, therefore $(V^{<\kappa}, T^{<\kappa})$ is closed under the basic set constructions. Furthermore if

$$\langle\langle a, b \rangle, \langle a', b' \rangle\rangle \in \llbracket \text{Fam} \rrbracket(V^{<\kappa}, T^{<\kappa}) ,$$

then

$$\langle \langle a, b \rangle, \langle a', b' \rangle \rangle \in \llbracket \text{Fam} \rrbracket (V^{<\beta}, T^{<\beta})$$

for some $\beta < \kappa$. Therefore there exists a $\rho < \kappa$ s.t.

$$\langle \langle f(a, b), g(a, b) \rangle, \langle f'(a', b'), g'(a', b') \rangle \rangle \in \llbracket \text{Fam} \rrbracket (V^{<\rho}, T^{<\rho}) \subseteq \llbracket \text{Fam} \rrbracket (V^{<\kappa}, T^{<\kappa}) .$$

This means that $V^{<\kappa}$ is a candidate for the interpretation of $\langle \widehat{U}_{f,g}, \widehat{U}_{f',g'} \rangle$, and therefore

$$\langle \widehat{U}_{f,g}, \widehat{U}_{f',g'} \rangle \in V^{\kappa+2} \subseteq V^{<M} .$$

That $\llbracket U_{f,g} \rrbracket_\rho$ has the closure properties needed in order to interpret $U_{f,g}$ follows by the construction.

The remaining steps are as for the model of simple $\text{ML}_W U$, i.e. we need $< \omega$ many admissibles in order to interpret the basic set constructions on top of V, T_V . Therefore the type theory can be interpreted in Kripke Platek set theory plus one recursively Mahlo ordinal and finitely many admissibles above it. So we have given the essence of a proof of the following theorem:

Theorem 5.3.

- (a) We can model $\text{ML}_W + (\text{Mahlo})$ in $(\text{KPM})^+$.
- (b) $|\text{ML}_W + (\text{Mahlo})| \leq |(\text{KPM})^+|$
- (c) The previous statements hold as well if we replace intensional by extensional equality.

By [30] we obtain the other direction and obtain therefore

Theorem 5.4. (a) $|\text{ML}_W + (\text{Mahlo})| = |(\text{KPM})^+|$

- (b) The previous statements hold as well if we replace intensional by extensional equality.

References

- [1] P. Aczel. The strength of Martin-Löf's intuitionistic type theory with one universe. In S. Miettinen and J. Väänänen, editors, *Proceedings of the symposium on Mathematical Logic (Oulu 1974)*. Report No. 2, Dept. of Philosophy, Univ of Helsinki, 1977.
- [2] T. Arai. Ordinal diagrams for Π_3 -reflection. *J. Symbolic Logic*, 65(3):1375 – 1394, 2000.
- [3] T. Arai. Ordinal diagrams for recursively Mahlo universes. *Arch. Math. Logic*, 39(5):353 – 391, 2000.
- [4] T. Arai. Proof theory for theories of ordinals. I. Recursively Mahlo ordinals. *Ann. Pure Appl. Logic*, 122(1 – 3):1 – 85, 2003.

- [5] T. Arai. Proof theory for theories of ordinals. II. Π_3 -reflection. *Ann. Pure Appl. Logic*, 129(1 – 3):39 – 92, 2004.
- [6] J. Barwise. *Admissible Sets and Structures. An Approach to Definability Theory*. Omega-series. Springer, 1975.
- [7] M. Benke, P. Dybjer, and P. Jansson. Universes for generic programs and proofs in dependent type theory. *Nordic J. of Computing*, 10(4):265–289, 2003.
- [8] P. Dybjer. A general formulation of simultaneous inductive-recursive definitions in type theory. *Journal of Symbolic Logic*, 65(2):525–549, June 2000.
- [9] P. Dybjer and A. Setzer. Finite axiomatizations of inductive and inductive-recursive definitions. In *Workshop on Generic Programming, Marstrand, Sweden, 18 June 1998*. http://www.cs.ruu.nl/people/johanj/programme_wgp98.html, 1998.
- [10] P. Dybjer and A. Setzer. A finite axiomatization of inductive-recursive definitions. In J.-Y. Girard, editor, *Typed Lambda Calculi and Applications*, volume 1581 of *Lecture Notes in Computer Science*, pages 129–146. Springer, April 1999.
- [11] P. Dybjer and A. Setzer. Indexed induction-recursion. In R. Kahle, P. Schroeder-Heister, and R. Stärk, editors, *Proof Theory in Computer Science*, volume 2183 of *Lecture Notes in Computer Science*, pages 93 – 113. Springer, 2001.
- [12] P. Dybjer and A. Setzer. Induction-recursion and initial algebras. *Annals of Pure and Applied Logic*, 124:1–47, 2003.
- [13] P. Dybjer and A. Setzer. Indexed induction-recursion. 64 pp. Submitted., 2005.
- [14] E. Griffor and M. Rathjen. The strength of some Martin-Löf type theories. *Arch. math. Logic*, 33:347 – 385, 1994.
- [15] M. Hofmann. Syntax and semantics of dependent types. In A. M. Pitts and P. Dybjer, editors, *Semantics and logics of computation*, pages 79 – 130, Cambridge, 1997. Cambridge University Press.
- [16] G. Jäger. *Theories for Admissible Sets: A Unifying Approach to Proof Theory*. Bibliopolis, Naples, 1986.
- [17] P. Martin-Löf. *Intuitionistic type theory*. Bibliopolis, Naples, 1984.
- [18] B. Nordström, K. Petersson, and J. Smith. *Programming in Martin-Löf’s type theory. An Introduction*. Oxford University-Press, Oxford, 1990. Out of print. Online version available via <http://www.cs.chalmers.se/Cs/Research/Logic/book/>.

- [19] B. Nordström, K. Petersson, and J. M. Smith. Martin-löf's type theory. In S. Abramsky, D. M. Gabbay, and T. S. E. Maibaum, editors, *Handbook of logic in computer science, Vol. 5*, pages 1 – 37. Oxford Univ. Press, 2000.
- [20] E. Palmgren. On universes in type theory. In G. Sambin and J. Smith, editors, *Twenty five years of constructive type theory*, pages 191 – 204, Oxford, 1998. Oxford University Press.
- [21] A. Ranta. *Type-theoretical Grammar*. Clarendon Press, Oxford, 1995.
- [22] M. Rathjen. Proof-theoretical analysis of KPM. *Arch. math. Logic*, 30:377 – 403, 1991.
- [23] M. Rathjen. Proof theory of reflection. *Ann. Pure Appl. Logic*, 68:181 – 224, 1994.
- [24] M. Rathjen. Recent advances in ordinal analysis: Π_2^1 -CA and related systems. *Bulletin of Symbolic Logic*, 1:468 – 485, 1995.
- [25] A. Setzer. *Proof theoretical strength of Martin-Löf Type Theory with W-type and one universe*. PhD thesis, Universität München, 1993. available from <http://www.cs.swan.ac.uk/~csetzer>.
- [26] A. Setzer. A model for a type theory with Mahlo Universe. 10pp. Preprint, available via [http://www.cs.swan.ac.uk/~articles/uppermahlo.pdf](http://www.cs.swan.ac.uk/~/articles/uppermahlo.pdf), 1996.
- [27] A. Setzer. An upper bound for the proof theoretical strength of Martin-Löf Type Theory with W-type and one Universe. 33pp. Draft, available via <http://www.cs.swan.ac.uk/~articles/1papdiss.dvi> or .ps or .pdf, 1996.
- [28] A. Setzer. An introduction to well-ordering proofs in Martin-Löf's type theory. In G. Sambin and J. Smith, editors, *Twenty-five years of constructive type theory*, pages 245 – 263, Oxford, 1998. Clarendon Press.
- [29] A. Setzer. Well-ordering proofs for Martin-Löf type theory. *Annals of Pure and Applied Logic*, 92:113 – 159, 1998.
- [30] A. Setzer. Extending Martin-Löf type theory by one Mahlo-universe. *Arch. Math. Log.*, 39:155 – 181, 2000.
- [31] A. Setzer. Proof theory of Martin-Löf type theory – an overview. *Mathematiques et Sciences Humaines*, 42 année, n°165:59 – 99, 2004.
- [32] A. Setzer. Universes in type theory part II – Π_3 -reflection. In preparation, 2005.
- [33] S. G. Simpson. *Subsystems of second-order arithmetic*. Springer, 1999.
- [34] A. Troelstra and D. v. Dalen. *Constructivism in Mathematics. An Introduction, Vol. II*. North-Holland, Amsterdam, 1988.