# A SIMPLIFIED PROOF OF THE EPSILON THEOREMS

STEFAN HETZL

**Abstract.** We formulate Hilbert's epsilon calculus in the context of expansion proofs. This leads to a simplified proof of the epsilon theorems by disposing of the need for prenexification, Skolemisation, and their respective inverse transformations. We observe that the natural notion of cut in the epsilon calculus is associative.

**§1. Introduction.** The epsilon calculus has been introduced by D. Hilbert as part of his programme for the foundation of mathematics. It is based on adding an operator $\varepsilon$ to first-order logic with the intended semantics that $\varepsilon_x A$ denotes an individual $x$ for which $A$ holds, if such an $x$ exists. The quantifiers of first-order logic can be defined by the $\varepsilon$-operator and thus first-order logic can be understood based on quantifier-free reasoning involving the $\varepsilon$-operator. The epsilon calculus is obtained from first-order logic by adding simple axioms fixing the behaviour of the $\varepsilon$-operator. The epsilon calculus has been used for consistency proofs [1], for studying proof complexity [2, 24], in computational proof theory [8], and in linguistics (see [6] for a survey).

Its main theorems, the epsilon theorems, are an analogue of cut-elimination in the sequent calculus. The first epsilon theorem states that, if a quantifier- and epsilon-free formula is provable in the epsilon calculus, then it is provable in first-order logic. The second epsilon theorem states that, if an epsilon-free formula is provable in the epsilon calculus, then it is provable in first-order logic. The standard proof of the epsilon theorems [17] (see also [24]) proceeds as follows: first, a (non-elementary) elimination procedure is carried out as the proof of the first epsilon theorem. For proving the second epsilon theorem, first, the extended first epsilon theorem is shown which states that if a formula $\exists \overline{x} A(\overline{x})$, with $A$ quantifier- and $\varepsilon$-free, is provable in the epsilon calculus, then there are term tuples $\overline{t_1}, \ldots, \overline{t_n}$ s.t. $\bigvee_{i=1}^{n} A(\overline{t_i})$ is provable in first-order logic. This result can be shown by, essentially, repeating the proof of the first epsilon theorem with some extra care concerning the disjunction. Finally, the second epsilon theorem is proved by Skolemisation, prenexification, and an application of the extended first epsilon theorem, followed by deskolemisation and deprenexification.

In this paper we give a simplified proof of the epsilon theorems. The crucial observation is that, if we phrase the epsilon theorems in the context of expansion proofs, then the proof of the first epsilon theorem already gives the second epsilon theorem as a result. Neither prenexification, nor Skolemisation, nor their inverses are required. Expansion trees are a formalism for representing Herbrand expansions of infix formulas. They have originally been introduced in the context of higher-order logic in [22]. Expansion trees have turned out to be useful for both theoretical investigations

as, e.g., in [7, 10, 26] as well as implementations (see, e.g., [3, 15, 19] and have become standard in the literature).

The second contribution of this paper is to develop the natural notion of cut in the epsilon calculus and to observe that it is associative, which suggests the perspective of investigating the epsilon calculus from the point of view of category theory along the lines of [18].

**§2. Terms, formulas, and substitutions.** To obtain a uniform treatment of first-order quantifiers and the $\varepsilon$-operator it is convenient to represent first-order terms and formulas, with and without $\varepsilon$, as simply typed lambda terms in the spirit of [11]. The *base types* of our simply typed lambda calculus are $\iota$, for individual, and $o$, for Boolean. The *types* are formed from the base types with the binary type constructor $\to$. The *logical constants* are $\wedge\colon o \to o \to o$, $\vee\colon o \to o \to o$, and, since we will only be dealing with first-order logic, $\forall\colon (\iota \to o) \to o$, $\exists\colon (\iota \to o) \to o$ and $\varepsilon\colon (\iota \to o) \to \iota$. A first-order function symbol is a (lambda) constant of type $\iota^n \to \iota$ for some $n \geq 0$. A first-order predicate symbol is a (lambda) constant of type $\iota^n \to o$ for some $n \geq 0$. A first-order language is a set of first-order function symbols and first-order predicate symbols. We will only consider first-order languages $L$ that are closed under dualisation, i.e., for every $P\colon \iota^n \to o \in L$ there is a *dual* $P^\perp\colon \iota^n \to o \in L$. *Lambda terms* are formed from variables, the logical constants, and the constants in $L$ using abstraction and application as usual. We identify terms which only differ in the names of bound variables; hence, $\alpha$-equivalence is denoted by $=$. We write $\mathrm{FV}(M)$ for the set of free variables of the term $M$. We employ the usual notational conventions such as the infix notation of binary connectives and the abbreviation of $\forall \lambda x.A$ as $\forall x A$, of $\exists \lambda x.A$ as $\exists x A$, and of $\varepsilon \lambda x.A$ as $\varepsilon_x A$.

A *first-order formula* is a lambda term $M$ of type $o$ s.t. all variables occurring in $M$ are of type $\iota$. Since, in this paper, we are dealing with first-order logic only, we will refer to a first-order formula simply as a *formula*. A *sentence* is a formula without free variables. A *first-order term* is a lambda term $M$ of type $\iota$ s.t. all variables occurring in $M$ are of type $\iota$ and $M$ does not contain $\forall$ nor $\exists$ (but may contain $\varepsilon$). An *$\varepsilon$-term* is a first-order term of the form $\varepsilon_x A$. Note that $\varepsilon$-terms do not contain quantifiers. A first-order term or formula is called *$\varepsilon$-free* if it does not contain $\varepsilon$. The $\varepsilon$-free terms and formulas of this paper are those of standard first-order logic. We define the dual of a formula inductively by $(P(t_1, \dots, t_n))^\perp = P^\perp(t_1, \dots, t_n)$ for an atom $P(t_1, \dots, t_n)$, $(A \wedge B)^\perp = A^\perp \vee B^\perp$, $(A \vee B)^\perp = A^\perp \wedge B^\perp$, $(\forall x\, A)^\perp = \exists x\, A^\perp$, and $(\exists x\, A)^\perp = \forall x\, A^\perp$. We sometimes abbreviate $A^\perp \vee B$ as $A \supset B$.

For the sake of precision, it will later be useful to indicate certain locations in formulas explicitly. To that aim, we define: a *position* is a finite word over the alphabet $\{1, 2\}$. We write $\langle\rangle$ for the empty position. If $q = q_1 \cdots q_n$ and $r = r_1 \cdots r_k$ are positions, their concatenation is $qr = q_1 \cdots q_n r_1 \cdots r_k$. For a set $R$ of positions we define $qR = \{qr \mid r \in R\}$. The positions of a formula are defined inductively as follows. $\mathrm{Pos}(A) = \{\langle\rangle\}$ if $A$ is an atom, $\mathrm{Pos}(A \circ B) = 1\mathrm{Pos}(A) \cup 2\mathrm{Pos}(B)$ for $\circ \in \{\wedge, \vee\}$, and $\mathrm{Pos}(QxA) = 1\mathrm{Pos}(A)$ for $Q \in \{\forall, \exists\}$.

Since we are dealing with first-order logic only, we limit our notion of substitution to replacing objects of type $\iota$ by first-order terms. On the other hand, in our setting, substitutions, in addition to replacing variables by first-order terms, may also replace closed $\varepsilon$-terms by first-order terms. Consequently we define: a *substitution*

is a finite set of pairs $\sigma = [x_1 \backslash M_1, \dots, x_m \backslash M_m, e_1 \backslash N_1, \dots, e_n \backslash N_n]$ where $x_1, \dots, x_m$ are pairwise different variables of type $\iota$, $e_1, \dots, e_n$ are pairwise different closed $\varepsilon$-terms and $M_1, \dots, M_m, N_1, \dots, N_n$ are first-order terms. The domain of $\sigma$ is $\mathrm{dom}(\sigma) = \{x_1, \dots, x_m, e_1, \dots, e_n\}$ and the range of $\sigma$ is $\mathrm{rng}(\sigma) = \bigcup_{i=1}^{m} \mathrm{FV}(M_i) \cup \bigcup_{i=1}^{n} \mathrm{FV}(N_i)$. The application of a substitution $\sigma$ to a term $M$ is defined, as usual, by induction on $M$ renaming bound variables if necessary. A consequence of this definition is that, if $e_1$ is a subterm of $e_2$, then $M[e_1 \backslash N_1, e_2 \backslash N_2] = M[e_2 \backslash N_2]$ for all $M$, $N_1$, and $N_2$. We define $\beta$-reduction as usual by $(\lambda x.M)N \rightarrow_\beta M[x \backslash N]$.

Even though the definition of such substitutions is straightforward, in applying them certain phenomena arise that warrant some additional caution. This is best illustrated by first observing that, if $x, y$ are variables of type $\iota$ and $s, t$ first-order terms with $y \notin \mathrm{FV}(t)$, then $[y \backslash s][x \backslash t] = [x \backslash t][y \backslash s[x \backslash t]]$. If, instead, $a$ is a closed $\varepsilon$-term, we no longer have $[y \backslash s][a \backslash t] = [a \backslash t][y \backslash s[a \backslash t]]$ as the following example shows.

EXAMPLE 1. *Let* $M = P(\varepsilon_x Q(x,y))$, $s = c$, $t = d$, *and* $a = \varepsilon_x Q(x,c)$, *then* $M[y \backslash s][a \backslash t] = P(a)[a \backslash d] = P(d)$ *but* $M[a \backslash t][y \backslash s[a \backslash t]] = M[y \backslash s] = P(a)$.

The crucial point in the above example is that the $\varepsilon$-term $a$ which is replaced by $d$ is formed partially from $M$ and partially from $s$ in $M[y \backslash s]$, a phenomenon that does not exist when substituting variables only. A way to obtain an analogous commutation property for substitution of $\varepsilon$-terms is, roughly speaking, to only replace such $\varepsilon$-terms whose structure is of maximal complexity among the terms considered. These $\varepsilon$-terms cannot be obtained by substitution from less complex $\varepsilon$-terms. In the epsilon calculus this is made precise by the notion of subordination and rank but since this is a more general phenomenon, we develop these notions for lambda terms here.

DEFINITION 2. *Let $M$ be a lambda term and let* x *be a variable, then the set of terms subordinate to $M$ w.r.t. x, $\mathrm{subord}_x(M)$, is defined inductively as follows*:

$$\mathrm{subord}_x(c) = \emptyset \text{ for a constant } c,$$
$$\mathrm{subord}_x(y) = \emptyset \text{ for a variable } y,$$
$$\mathrm{subord}_x(MN) = \mathrm{subord}_x(M) \cup \mathrm{subord}_x(N),$$
$$\mathrm{subord}_x(\lambda y.M) = \begin{cases} \{\lambda y.M\}, & \text{if } x \in \mathrm{FV}(\lambda y.M), \\ \emptyset, & \text{otherwise}. \end{cases}$$

Let $\lambda x.M$ and $N$ be lambda terms, then $N$ is called *subordinate* to $\lambda x.M$ if $N \in \mathrm{subord}_x(M)$, in other words, if $N$ is a subterm of $M$ that starts with an abstraction and contains $x$ as free variable.

DEFINITION 3. *The* rank *of a lambda term is defined inductively by* $\mathrm{rk}(x) = \mathrm{rk}(c) = 0$, $\mathrm{rk}(MN) = \max\{\mathrm{rk}(M), \mathrm{rk}(N)\}$, *and* $\mathrm{rk}(\lambda x.M) = 1 + \max\{\mathrm{rk}(N) \mid N \in \mathrm{subord}_x(M)\}$ *with the convention that* $\max \emptyset = 0$.

EXAMPLE 4. *Let* $N = \lambda y.yx$ *and* $M = \lambda x.zN$, *then $N$ is subordinate to $M$,* $\mathrm{rk}(N) = 1$ *and* $\mathrm{rk}(M) = 2$.

Note that this generalises the traditional definition of rank of an $\varepsilon$-term [17, 24]. Also note that this formalism, in contrast to the traditional one [17, 24], neither requires notions such as semi-terms and semi-formulas nor does it require the variables bound by different $\varepsilon$'s to be different, which, although it is theoretically trivial, is a nuisance for implementations.

LEMMA 5. *For any lambda term $M$ and any substitution $\sigma$*: $\mathrm{rk}(M\sigma) = \mathrm{rk}(M)$.

*Proof.* First observe that, for any lambda term $M$, $x \notin \mathrm{FV}(M)$ implies $\mathrm{subord}_x(M) = \emptyset$. Secondly, we claim that

$$\mathrm{subord}_x(M\sigma) = \mathrm{subord}_x(M)\sigma \qquad (*)$$

for any lambda term $M$, any variable $x$, and any substitution $\sigma$ with $x \notin \mathrm{dom}(\sigma) \cup \mathrm{rng}(\sigma)$. This is shown by induction on the structure of $M$ with the cases for constants and application being straightforward and the one for variables relying on $x \notin \mathrm{rng}(\sigma)$. For abstraction let $M = \lambda y.N$. If $x \notin \mathrm{FV}(M)$ we are done by the above observation, so assume that $x \in \mathrm{FV}(M)$. Then $\mathrm{subord}_x(\lambda y.N)\sigma = \{\lambda y.N\}\sigma = \{\lambda y'.N[y\backslash y']\sigma\}$ where $y' \notin \mathrm{dom}(\sigma) \cup \mathrm{rng}(\sigma) \cup \mathrm{FV}(M)$. On the other hand, we also have $x \in \mathrm{FV}(N[y\backslash y']\sigma)$ because $x \notin \mathrm{dom}(\sigma)$ and therefore $\mathrm{subord}_x((\lambda y.N)\sigma) = \mathrm{subord}_x(\lambda y'.N[y\backslash y']\sigma) = \{\lambda y'.N[y\backslash y']\sigma\}$.

For showing the lemma it suffices to consider the case $M = \lambda x.M_0$. We proceed by induction on the rank of $M$. We assume w.l.o.g. that $x \notin \mathrm{dom}(\sigma) \cup \mathrm{rng}(\sigma)$. Then we have

$$\begin{aligned}
\mathrm{rk}(M\sigma) &= 1 + \max\{\mathrm{rk}(N) \mid N \in \mathrm{subord}_x(M\sigma)\} =^{(*)} 1 + \max\{\mathrm{rk}(N) \mid N \in \mathrm{subord}_x(M)\sigma\} \\
&= 1 + \max\{\mathrm{rk}(N_0\sigma) \mid N_0 \in \mathrm{subord}_x(M)\} = 1 + \max\{\mathrm{rk}(N_0) \mid N_0 \in \mathrm{subord}_x(M)\} \\
&= \mathrm{rk}(M). \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \square
\end{aligned}$$

LEMMA 6. *Let $y$ be a variable, let $a$ be a closed $\varepsilon$-term, let $s$ and $t$ be first-order terms s.t. $y \notin \mathrm{FV}(t)$, and let $M$ be a lambda term s.t. $\mathrm{rk}(\lambda y.M) \leq \mathrm{rk}(a)$. Then $M[y\backslash s][a\backslash t] = M[a\backslash t][y\backslash s[a\backslash t]]$.*

*Proof.* Suppose $M$ contains an $\varepsilon$-term $a' = \lambda x.A$ s.t. $a'[y\backslash s] = a$, then $\lambda x.A \in \mathrm{subord}_y(M)$ and hence $\mathrm{rk}(\lambda y.M) > \mathrm{rk}(a')$. On the other hand, by Lemma 5, we have $\mathrm{rk}(a') = \mathrm{rk}(a)$, contradiction. $\qquad\square$

Thus, when replacing $\varepsilon$-terms of maximal rank, the substitution has the above commutation property that is well-known from ordinary substitutions (which only replace variables).

**§3. Expansion proofs.** In this section we describe (a variant of) expansion trees, originally introduced in [22] in the setting of higher-order logic, that is suitable for our purposes. In particular witness terms may contain $\varepsilon$'s. It will be convenient to treat universal quantifiers with Skolem symbols instead of eigenvariables, similarly to the Skolem expansion trees of [22]. Therefore our expansion trees will not contain free variables.

Skolemisation is often applied in a refutational setting for obtaining a satisfiability-equivalent formula by replacing existential quantifiers by new function symbols. For expansion proofs, the dual transformation that yields a validity-equivalent formula by replacing universal quantifiers is needed. In this validity-preserving Skolemisation, sometimes also called Herbrandisation, of a formula $A$, an occurrence of a quantified subformula $\forall y\, C$ is replaced by $C[y\backslash f(x_1, \ldots, x_n)]$ where $x_1, \ldots, x_n$ are the variables bound by existential quantifiers on the path from the root of $A$ to the root of $C$. In a detailed inductive definition of expansion trees we have to also consider the case of intermediate subformula occurrences, i.e., of $B$ s.t. $A$ is $A[B[\forall y\, C]]$. In $B$ the variables

$x_1, \dots, x_k$, for some $k \leq n$, are already free and thus are considered parameters of $B$ while the quantifiers $\exists x_{k+1}, \dots, \exists x_n$ are in $B$. To deal with this situation precisely we introduce the notion of Skolem mapping.

DEFINITION 7. *Let $A$ be a sentence. A* Skolem mapping *for A is a pair $(\overline{p}; s)$ s.t. 1. $\overline{p}$ is a* k-*tuple of terms without free variables for some $k \geq 0$ and* 2. s *assigns a function symbol $s_q$ to every position* q *of a universal quantifier in A s.t.:*

3. *$s_q$ does not occur in $A$,*
4. *$q_1 \neq q_2$ implies $s_{q_1} \neq s_{q_2}$, and*
5. *the arity of $s_q$ is* k *plus the number of existential quantifiers on the path from the root of $A$ to* q.

The tuple $\overline{p}$ are the *parameters* of $(\overline{p}, s)$. The *range* of $(\overline{p}, s)$ is the range of $s$.

EXAMPLE 8. *$\exists x\,(P^\perp(x) \vee P(f(x)))$ is the Skolemisation of $\exists x \forall y\,(P^\perp(x) \vee P(y))$; $((); \{1 \mapsto f\})$ is a Skolem mapping for $\exists x \forall y (P^\perp(x) \vee P(y))$. In a proof of $\exists x \forall y\,(P^\perp(x) \vee P(y))$, the variable* x *may be instantiated with some term* t. *This will be reflected by using the Skolem mapping $((t); \{\langle\rangle \mapsto f\})$ for the formula $\forall y\,(P^\perp(t) \vee P(y))$ in order to store the information that* y *should be instantiated with $f(t)$.*

DEFINITION 9. *We define the notion of* expansion tree w.r.t. a Skolem mapping *simultaneously with its* shallow mapping $\mathrm{Sh}(\cdot)$ *and* deep mapping $\mathrm{Dp}(\cdot)$.

1. *Every atom $A$ without free variables is an expansion tree w.r.t. any Skolem mapping $(\overline{t}; \emptyset)$. We define $\mathrm{Sh}(A) = A$ and $\mathrm{Dp}(A) = A$.*
2. *If $E_1$ and $E_2$ are expansion trees w.r.t. Skolem mappings $(\overline{t}; s_1)$ and $(\overline{t}; s_2)$ respectively and $\circ \in \{\wedge, \vee\}$, then $E_1 \circ E_2$ is an expansion tree w.r.t. $(\overline{t}; 1s_1 \cup 2s_2)$. We define $\mathrm{Sh}(E_1 \circ E_2) = \mathrm{Sh}(E_1) \circ \mathrm{Sh}(E_2)$ and $\mathrm{Dp}(E_1 \circ E_2) = \mathrm{Dp}(E_1) \circ \mathrm{Dp}(E_2)$.*
3. *If $\{t_1, \dots, t_n\}$ is a set of terms and $E_1, \dots, E_n$ are expansion trees w.r.t. the Skolem mappings $(\overline{p}, t_1; s), \dots, (\overline{p}, t_n; s)$ and $\mathrm{Sh}(E_i) = A[x \backslash t_i]$ for $i = 1, \dots, n$, then $E = \exists x\, A +^{t_1} E_1 \cdots +^{t_n} E_n$ is an expansion tree w.r.t. the Skolem mapping $(\overline{p}; 1s)$. We define $\mathrm{Sh}(E) = \exists x\, A$ and $\mathrm{Dp}(E) = \bigvee_{i=1}^{n} \mathrm{Dp}(E_i)$.*
4. *If $E$ is an expansion tree w.r.t. a Skolem mapping $(\overline{p}, s)$ with $\mathrm{Sh}(E) = A[x \backslash f(\overline{p})]$ where* f *does not occur in $A$, then $E' = \forall x\, A +^{f(\overline{p})} E$ is an expansion tree with Skolem mapping $(\overline{p}, 1s \cup \{\langle\rangle \mapsto f\})$, $\mathrm{Sh}(E') = \forall x\, A$ and $\mathrm{Dp}(E') = \mathrm{Dp}(E)$.*

*We simply say that $E$ is an* expansion tree *if there exists a Skolem mapping w.r.t. which $E$ is an expansion tree.*

Note that $\varepsilon$, being part of our logic, may appear anywhere in an expansion tree, subject to the general restriction that an $\varepsilon$-term is quantifier-free. If we want to consider formulas, expansion trees, etc. that do not contain $\varepsilon$ we explicitly designate them as *$\varepsilon$-free* formulas, expansion trees, etc.

A *sequent* is a finite set of sentences. If $\Gamma = A_1, \dots, A_n$ is a sequent, $E_1, \dots, E_n$ are expansion trees with $\mathrm{Sh}(E_i) = A_i$ for $i = 1, \dots, n$ and Skolem mappings of disjoint range, then $S = E_1, \dots, E_n$ is called *expansion sequent* with $\mathrm{Sh}(S) = \Gamma$ and $\mathrm{Dp}(S) = \bigvee_{i=1}^{n} \mathrm{Dp}(E_i)$.

DEFINITION 10. *An expansion sequent $S$ with $\mathrm{Sh}(S) = \Gamma$ is called* expansion proof of $\Gamma$ *if $\mathrm{Dp}(S)$ is a tautology.*

EXAMPLE 11. *Let* s *and* t *be closed first-order terms, then*

$$E_0(s, t) = \forall y \, (P^\perp(s) \vee P(y)) +^{f(t)} P^\perp(s) \vee P(f(t))$$

*is an expansion tree of* $\forall y \, (P^\perp(s) \vee P(y))$ *with Skolem mapping* $((t); \{\langle\rangle \mapsto f\})$. *Let* $E(t) = E_0(t, t)$, *then*

$$D_0 = \exists x \forall y \, (P^\perp(x) \vee P(y)) +^c E(c) \; and$$
$$D = \exists x \forall y (P^\perp(x) \vee P(y)) +^c E(c) +^{f(c)} E(f(c))$$

*are expansion trees of* $\exists x \forall y \, (P^\perp(x) \vee P(y))$, *both with Skolem mapping* $((); \{1 \mapsto f\})$. $D$ *is an expansion proof but* $D_0$ *is not.*

In the setting of an existential formula $\exists x \, A$ where $A$ is quantifier-free, two finite sets of instances, $\bigvee_{t \in T} A[x \backslash t]$ and $\bigvee_{u \in U} A[x \backslash u]$, can be merged by simply forming their union $\bigvee_{t \in T \cup U} A[x \backslash t]$. We will now generalise this merge operation to expansion trees. Two expansion trees $E_1$ and $E_2$ are called *mergeable* if $\mathrm{Sh}(E_1) = \mathrm{Sh}(E_2)$ and they have the same Skolem mapping. If two expansion trees $E_1$ and $E_2$ satisfy $\mathrm{Sh}(E_1) = \mathrm{Sh}(E_2)$ we can w.l.o.g. assume that they are mergeable by renaming the Skolem symbols in one of them.

DEFINITION 12. *Let* $E_1$ *and* $E_2$ *be mergeable expansion trees, then their* merge $E_1 \sqcup E_2$ *is defined as follows*:

1. *If* $A$ *is an atom, then* $E_1 = E_2$ *and we define*

$$E_1 \sqcup E_2 = E_1 = E_2.$$

2. *If* $A = A' \circ A''$ *for* $\circ \in \{\wedge, \vee\}$, *then* $E_1 = E_1' \circ E_1''$, $E_2 = E_2' \circ E_2''$ *and we define*

$$E_1 \sqcup E_2 = (E_1' \sqcup E_2') \circ (E_1'' \sqcup E_2'').$$

3. *If* $A = \exists x \, A'$, *then* $E_1 = \exists x \, A' +^{r_1} E_{1,1} \ldots +^{r_k} E_{1,k} +^{s_1} F_1 \ldots +^{s_l} F_l$ *and* $E_2 = \exists x \, A' +^{r_1} E_{2,1} \ldots +^{r_k} E_{2,k} +^{t_1} G_1 \ldots +^{t_m} G_m$ *where* $\{s_1, \ldots, s_l\} \cap \{t_1, \ldots, t_m\} = \emptyset$ *and we define*

$$E_1 \sqcup E_2 = \exists x \, A' +^{r_1} (E_{1,1} \sqcup E_{2,1}) \ldots +^{r_k} (E_{1,k} \sqcup E_{2,k})$$
$$+^{s_1} F_1 \ldots +^{s_l} F_l +^{t_1} G_1 \ldots +^{t_m} G_m.$$

4. *If* $A = \forall x \, A'$, *then* $E_1 = \forall x \, A' +^{s\langle\rangle^{(\overline{p})}} E_1'$, $E_2 = \forall x \, A' +^{s\langle\rangle^{(\overline{p})}} E_2'$ *and we define*

$$E_1 \sqcup E_2 = \forall x \, A' +^{s\langle\rangle^{(\overline{p})}} (E_1' \sqcup E_2').$$

Two expansion sequents $S = E_1, \ldots, E_n$ and $T = F_1, \ldots, F_m$ are called mergeable if $\mathrm{Sh}(E_i) = \mathrm{Sh}(F_j)$ implies that $E_i$ and $F_j$ are mergeable for all $i \in \{1, \ldots, n\}$ and $j \in \{1, \ldots, m\}$. Note that, up to renaming Skolem symbols, any two expansion sequents are mergeable. Let $S_1 = E_{1,1}, \ldots, E_{1,k}, F_1, \ldots, F_l$ and $S_2 = E_{2,1}, \ldots, E_{2,k}, G_1, \ldots, G_m$ be mergeable expansion sequents s.t. $\mathrm{Sh}(E_{1,i}) = \mathrm{Sh}(E_{2,i})$ for $i = 1, \ldots, k$ and $\mathrm{Sh}(F_1, \ldots, F_n) \cap \mathrm{Sh}(G_1, \ldots, G_m) = \emptyset$. Then their merge is defined as $S_1 \sqcup S_2 = E_{1,1} \sqcup E_{2,1}, \ldots, E_{1,k} \sqcup E_{2,k}, F_1, \ldots, F_l, G_1, \ldots G_m$. Note that $\mathrm{Sh}(E_1 \sqcup E_2) = \mathrm{Sh}(E_1) = \mathrm{Sh}(E_2)$ and $\mathrm{Sh}(S_1 \sqcup S_2) = \mathrm{Sh}(S_1) \cup \mathrm{Sh}(S_2)$. Moreover, note that the merge operation, both on expansion trees and on expansion sequents, is associative. We now turn to analysing the deep mapping of a merge. For quantifier-free sentences $A, B$ we write $A \Rightarrow B$ if the formula $A \supset B$ is a tautology.

LEMMA 13. *Let $E_1, E_2$ be mergeable expansion trees, then $\mathrm{Dp}(E_1) \vee \mathrm{Dp}(E_2) \Rightarrow \mathrm{Dp}(E_1 \sqcup E_2)$. Let $S_1, S_2$ be mergeable expansion sequents, then $\mathrm{Dp}(S_1) \vee \mathrm{Dp}(S_2) \Rightarrow \mathrm{Dp}(S_1 \sqcup S_2)$.*

*Proof.* The result on formulas is proved by a straightforward induction on the structure of $\mathrm{Sh}(E_1) = \mathrm{Sh}(E_2)$. The most interesting case is that of $\wedge$ since it hinders the logical equivalence: For $E_1 = E_1' \wedge E_1''$ and $E_2 = E_2' \wedge E_2''$ we have

$$
\begin{aligned}
\mathrm{Dp}(E_1) \vee \mathrm{Dp}(E_2) &= (\mathrm{Dp}(E_1') \wedge \mathrm{Dp}(E_1'')) \vee (\mathrm{Dp}(E_2') \wedge \mathrm{Dp}(E_2'')) \\
&\Rightarrow (\mathrm{Dp}(E_1') \vee \mathrm{Dp}(E_2')) \wedge (\mathrm{Dp}(E_1'') \vee \mathrm{Dp}(E_2'')) \\
&\Rightarrow^{\mathrm{IH}} \mathrm{Dp}(E_1' \sqcup E_2') \wedge \mathrm{Dp}(E_1'' \sqcup E_2'') \\
&= \mathrm{Dp}(E_1 \sqcup E_2).
\end{aligned}
$$

The result on sequents follows from that on formulas since we have

$$
\begin{aligned}
\mathrm{Dp}(S_1) \vee \mathrm{Dp}(S_2) &= \bigvee_{i=1}^{k} (\mathrm{Dp}(E_{1,i}) \vee \mathrm{Dp}(E_{2,i})) \bigvee_{i=1}^{l} \mathrm{Dp}(F_i) \bigvee_{i=1}^{m} \mathrm{Dp}(G_i) \\
&\Rightarrow \bigvee_{i=1}^{k} \mathrm{Dp}(E_{1,i} \sqcup E_{2,i}) \bigvee_{i=1}^{l} \mathrm{Dp}(F_i) \bigvee_{i=1}^{m} \mathrm{Dp}(G_i) \\
&= \mathrm{Dp}(S_1 \sqcup S_2). \qquad \square
\end{aligned}
$$

DEFINITION 14. *Let $E$ be an expansion tree and let $\sigma$ be a substitution, then $E\sigma$ is defined as follows*:

1. *If $E$ is an atom, then $E\sigma$ is already defined.*
2. *If $E = E_1 \circ E_2$ for $\circ \in \{\wedge, \vee\}$, then $E\sigma = E_1\sigma \circ E_2\sigma$.*
3. *If $E = \exists x\, A +^{t_1} E_1 \cdots +^{t_n} E_n$, let $\{s_1, \dots s_k\} = \{t_1\sigma, \dots, t_n\sigma\}$ and define*

$$
E\sigma = \exists x\, A\sigma +^{s_1} \bigsqcup_{\substack{1 \le i \le n \\ t_i\sigma = s_1}} E_i\sigma \cdots +^{s_k} \bigsqcup_{\substack{1 \le i \le n \\ t_i\sigma = s_k}} E_i\sigma.
$$

   *This is well-defined since $t_i\sigma = t_j\sigma$ implies that $E_i\sigma$ and $E_j\sigma$ are mergeable.*
4. *If $E = \forall x\, A +^{f(\bar{t})} E_0$, then $E\sigma = \forall x\, A\sigma +^{f(\bar{t})\sigma} E_0\sigma$.*

*For an expansion sequent $S = E_1, \dots, E_n$ we define $S\sigma$ as the merge $\bigsqcup_{i=1}^{n} E_i\sigma$ of the singleton expansion sequents $E_1\sigma, \dots, E_n\sigma$.*

We have $\mathrm{Sh}(E\sigma) = \mathrm{Sh}(E)\sigma$ and $\mathrm{Sh}(S\sigma) = \mathrm{Sh}(S)\sigma$. The deep formula has the following behaviour under substitution:

LEMMA 15. *Let $E$ be an expansion tree, let $S$ be an expansion sequent, and let $\sigma$ be a substitution. Then $\mathrm{Dp}(E)\sigma \Rightarrow \mathrm{Dp}(E\sigma)$ and $\mathrm{Dp}(S)\sigma \Rightarrow \mathrm{Dp}(S\sigma)$.*

*Proof.* For the first statement we proceed by induction on $E$. All cases except the existential quantifier are straightforward. For the existential quantifier, following the notation of Definition 14, we have $\mathrm{Dp}(\exists x\, A +^{t_1} E_1 \cdots +^{t_n} E_n)\sigma = \bigvee_{i=1}^{n} \mathrm{Dp}(E_i)\sigma \Rightarrow^{\mathrm{IH}} \bigvee_{i=1}^{n} \mathrm{Dp}(E_i\sigma) = \bigvee_{j=1}^{k} \bigvee_{\substack{1 \le i \le n \\ t_i\sigma = s_j}} \mathrm{Dp}(E_i\sigma) \Rightarrow^{\mathrm{Lem.\ 13}} \bigvee_{j=1}^{k} \mathrm{Dp}(\bigsqcup_{\substack{1 \le i \le n \\ t_i\sigma = s_j}} E_i\sigma)$.

For the second statement let $S = E_1, \dots, E_n$ and observe that $\mathrm{Dp}(S)\sigma = (\bigvee_{i=1}^{n} \mathrm{Dp}(E_i))\sigma = \bigvee_{i=1}^{n} \mathrm{Dp}(E_i)\sigma \Rightarrow \bigvee_{i=1}^{n} \mathrm{Dp}(E_i\sigma) \Rightarrow^{\mathrm{Lem.\ 13}} \mathrm{Dp}(\bigsqcup_{i=1}^{n} E_i\sigma) = \mathrm{Dp}(S\sigma)$. $\square$

In particular, applying a substitution to an expansion proof yields an expansion proof.

LEMMA 16. *If an $\varepsilon$-free sequent $\Gamma$ has an expansion proof, then $\Gamma$ has an $\varepsilon$-free expansion proof.*

*Proof.* Let $S$ be an expansion proof of $\Gamma$, then neither $S$ nor $\Gamma$ contains a free variable and thus all $\varepsilon$-terms in $S$ are closed. Let $e_1, \ldots, e_n$ be the $\varepsilon$-terms in $S$ and let $t$ be any $\varepsilon$-free term. Then $S' = S[e_1 \backslash t, \ldots, e_n \backslash t]$ is an $\varepsilon$-free expansion proof of $\Gamma[e_1 \backslash t, \ldots, e_n \backslash t] = \Gamma$.                      □

It is straightforward to translate a cut-free sequent calculus proof into an expansion proof and vice versa. These translations are described in detail in [7], here we just recall the essential points: the most natural match for the expansion proofs used in this paper is a cut-free sequent calculus where the rule for the universal quantifier is

$$\frac{\Gamma, A[x \backslash f(t_1, \ldots, t_n)]}{\Gamma, \forall x \, A} \ \forall,$$

where $t_1, \ldots, t_n$ are the terms inserted for the weak quantifiers on the path from $\forall x \, A$ to the end-sequent and a global condition on Skolem symbols guarantees their consistent choice.

An expansion proof is then read off from a cut-free sequent calculus proof recursively. The merge operation is used on (explicit and implicit) contractions. The joint trace of formula successors from the conclusion of the contraction to the end-sequent ensures mergeability. The Skolem condition on the sequent calculus proof ensures the existence of a Skolem mapping for the expansion proof. The invariant of the recursive extraction algorithm is that $\mathrm{Dp}(\Gamma)$ is a tautology where $\Gamma$ is the current end-sequent.

The translation in the other direction is a bottom-up construction of a sequent calculus proof guided by the given expansion proof. At each step, an outermost node of the expansion proof is removed by translating it to an inference in the sequent calculus. The use of invertible rules ensures that $\mathrm{Dp}(\Gamma)$ remains a tautology where $\Gamma$ is the current end-sequent. In case the universal quantifiers are modelled using eigenvariables, an acyclicity condition which mimics the acyclicity of the subterm order on Skolem terms is satisfied (see [7] for details).

**§4. The epsilon theorems.**    In this section we define epsilon proofs and show an elimination theorem which yields the epsilon theorems as corollaries.

DEFINITION 17.  *A critical axiom is a sentence of the form $A[x \backslash t] \supset A[x \backslash \varepsilon_x A]$.*

Note that, given a critical axiom $A' \supset A''$, the terms $t$ and $\varepsilon_x A$ are uniquely determined. Given a critical axiom $A' \supset A''$ we can therefore speak about *its $\varepsilon$-term $\varepsilon_x A$* and thus, given a set $C$ of critical axioms and an $\varepsilon$-term $\varepsilon_x A$ we can speak about the subset $C' \subseteq C$ of *critical axioms of $\varepsilon_x A$ in $C$*.

DEFINITION 18.  *An $\varepsilon$-preproof is a pair $P = (C; S)$ where $C$ is a finite set of critical axioms and $S$ is an expansion sequent. An $\varepsilon$-preproof $(C; S)$ is called $\varepsilon$-proof if $\bigwedge C \supset \mathrm{Dp}(S)$ is a tautology.*

Note that the notion of $\varepsilon$-proof is a straightforward generalisation of the notion of expansion proof which corresponds to the case $C = \emptyset$. In Section 5 we will show how

to read off $\varepsilon$-proofs from sequent calculus proofs. An $\varepsilon$-term $a$ is called *critical* in an epsilon preproof $P = (C; S)$ if there is a critical axiom of $a$ in $C$.

DEFINITION 19. *Let $P = (C; S)$ be an $\varepsilon$-preproof, let $a = \varepsilon_x A$ be a critical $\varepsilon$-term of maximal rank in $P$, and let $C_a = \{A[x \backslash t_i] \supset A[x \backslash a] \mid 1 \leq i \leq n\} \subseteq C$ be the set of critical formulas of $a$. The* reduct *of P w.r.t. a in C is defined as $P' = (C', S')$ where*

$$C' = (C \setminus C_a) \cup (C \setminus C_a)[a \backslash t_1] \cup \cdots \cup (C \setminus C_a)[a \backslash t_n] \text{ and}$$
$$S' = S \sqcup S[a \backslash t_1] \sqcup \cdots \sqcup S[a \backslash t_n].$$

*If $P'$ is the reduct of $P$ w.r.t. a we write $P \to^a P'$. If we want to ignore the $\varepsilon$-term a we write $P \to P'$. We write $\twoheadrightarrow$ for the reflexive and transitive closure of $\to$.*

For $P'$ to be an $\varepsilon$-proof we need to ensure in particular that $C'$ is a set of critical axioms. This will be obtained from the assumption that $\mathrm{rk}(a)$ is maximal via the following lemma.

LEMMA 20. *Let $a = \varepsilon_x A$ and $b = \varepsilon_y B$ be closed $\varepsilon$-terms, let $B[y \backslash s] \supset B[y \backslash b]$ be a critical axiom s.t. $\mathrm{rk}(a) \geq \mathrm{rk}(b)$, then $(B[y \backslash s] \supset B[y \backslash b])[a \backslash t]$ is a critical axiom of $b[a \backslash t]$.*

*Proof.* By Lemma 6 we have $(B[y \backslash s] \supset B[y \backslash b])[a \backslash t] = B'[y \backslash s[a \backslash t]] \supset B'[y \backslash b [a \backslash t]]$ where $B' = B[a \backslash t]$. This is a critical axiom of $b[a \backslash t]$. $\square$

LEMMA 21. *If P is an $\varepsilon$-proof and $P \to P'$, then $P'$ is an $\varepsilon$-proof.*

*Proof.* Let $P = (C; S), P' = (C'; S'), P \to^a P'$ where $a = \varepsilon_x A$ and $C_a = \{A[x \backslash t_i] \to A[x \backslash a] \mid 1 \leq i \leq n\}$. Then $C'$ and $S'$ are as in Definition 19. By Lemma 20, $C'$ consists of critical axioms; hence, $P'$ is a $\varepsilon$-preproof. Let $i \in \{1, \ldots, n\}$. Note that $a = \varepsilon_x A$ does not occur in $A$ and hence $A[x \backslash a][a \backslash t_i] = A[x \backslash t_i]$. So $C_a[a \backslash t_i] = \{B_j \to A[x \backslash t_i] \mid 1 \leq j \leq n\}$ for some formulas $B_1, \ldots, B_n$. Therefore $A[x \backslash t_i] \supset \bigwedge C_a[a \backslash t_i]$ is a tautology. Since $\bigwedge C \supset \mathrm{Dp}(S)$ is a tautology, so is $(\bigwedge C \supset \mathrm{Dp}(S))[a \backslash t_i]$ and, by Lemma 15, also $\bigwedge C[a \backslash t_i] \supset \mathrm{Dp}(S[a \backslash t_i])$. Therefore

$$A[x \backslash t_i] \supset \underbrace{\bigwedge (C \setminus C_a)[a \backslash t_i] \supset \mathrm{Dp}(S[a \backslash t_i])}_{\psi_i}$$

is a tautology. Moreover, $\bigwedge_{i=1}^{n} \neg A[x \backslash t_i] \supset \bigwedge C_a$ is a tautology and therefore

$$\bigwedge_{i=1}^{n} \neg A[x \backslash t_i] \supset \underbrace{\bigwedge (C \setminus C_a) \supset \mathrm{Dp}(S)}_{\psi_0}$$

is a tautology. Removing the case distinction shows that $\psi_0 \vee \psi_1 \vee \cdots \vee \psi_n$ is a tautology and thus, so is

$$\bigwedge C' \to \mathrm{Dp}(S) \vee \mathrm{Dp}(S[a \backslash t_1]) \vee \cdots \vee \mathrm{Dp}(S[a \backslash t_n]).$$

Therefore, by Lemma 13, also

$$\bigwedge C' \to \mathrm{Dp}(S')$$

is a tautology. $\square$

THEOREM 22. *Let $\Gamma$ be an $\varepsilon$-free sequent and let $(C; S)$ be an $\varepsilon$-proof of $\Gamma$. Then there is an expansion proof $S^*$ of $\Gamma$ s.t. $(C, S) \twoheadrightarrow (\emptyset, S^*)$.*

*Proof.* The *order* of a expansion proof $P$ w.r.t. $r \in \mathbb{N}$ is the number of $\varepsilon$-terms of rank $r$ in $P$. At each step of the reduction we pick an $\varepsilon$-term $a$ of maximal rank s.t. if $b$ is another $\varepsilon$-term of maximal rank, then $a$ is not a subterm of $b$. Thus substituting for $a$ does not change the other $\varepsilon$-terms of maximal rank. Hence the order w.r.t. the maximal rank strictly decreases and, once the last $\varepsilon$-term of maximal rank is eliminated, the maximal rank strictly decreases.     □

COROLLARY 23 (First $\varepsilon$-theorem). *If a quantifier- and $\varepsilon$-free sequent $\Gamma$ has an $\varepsilon$-proof, then $\Gamma$ is a tautology.*

*Proof.* By Theorem 22 there is an expansion proof $S^*$ of $\Gamma$. Since $\Gamma$ does not contain quantifiers, $\mathrm{Dp}(S^*) = \Gamma$ and thus $\Gamma$ is a tautology.     □

COROLLARY 24 (Extended first $\varepsilon$-theorem). *If a formula $\exists \overline{x}\, \varphi$, where $\varphi$ is quantifier- and $\varepsilon$-free, has an $\varepsilon$-proof then there are $\varepsilon$-free term tuples $\overline{t_1}, \dots, \overline{t_k}$ s.t. $\bigvee_{i=1}^{k} \varphi[\overline{x} \backslash \overline{t_i}]$ is a tautology.*

*Proof.* By Theorem 22 there is an expansion proof $S^*$ of $\exists \overline{x}\, \varphi$. By Lemma 16 we can assume that $S^*$ is $\varepsilon$-free. Since $\varphi$ does not contain quantifiers, $\mathrm{Dp}(S^*)$ is a formula of the required form.     □

COROLLARY 25 (Second $\varepsilon$-theorem). *If an $\varepsilon$-free sequent $\Gamma$ has an $\varepsilon$-proof, then $\Gamma$ has an $\varepsilon$-free expansion proof.*

*Proof.* Immediate from Theorem 22 and Lemma 16.     □

**§5. Cut for epsilon proofs.**  In this section we define a cut operation for $\varepsilon$-proofs which yields an $\varepsilon$-proof of $\Gamma, \Delta$ from $\varepsilon$-proofs of $\Gamma, A$ and $A^\perp, \Delta$. In particular, this operation allows to read off an $\varepsilon$-proof from a sequent calculus proof with cut in a straightforward way. We first recall the $\varepsilon$-translation of formulas from [17] (see also [24]).

DEFINITION 26. *The epsilon translation of formulas and terms is defined inductively as follows*:

$$c^\varepsilon = c, \qquad x^\varepsilon = x, \qquad f(t_1, \dots, t_n)^\varepsilon = f(t_1^\varepsilon, \dots, t_n^\varepsilon),$$
$$P(t_1, \dots, t_n)^\varepsilon = P(t_1^\varepsilon, \dots, t_n^\varepsilon), \quad (A \wedge B)^\varepsilon = A^\varepsilon \wedge B^\varepsilon, \qquad (A \vee B)^\varepsilon = A^\varepsilon \vee B^\varepsilon,$$
$$(\exists x\, A)^\varepsilon = A^\varepsilon[x \backslash \varepsilon_x A^\varepsilon], \qquad (\forall x\, A) = A^\varepsilon[x \backslash \varepsilon_x A^{\perp \varepsilon}], \qquad (\varepsilon_x A)^\varepsilon = \varepsilon_x A^\varepsilon.$$

Note that $A^{\perp \varepsilon} = A^{\varepsilon \perp}$ for all formulas $A$ which can be shown by a straightforward induction. Moreover, note that $A^\varepsilon = A$ for a quantifier-free formula $A$. In particular, $A^{\varepsilon \varepsilon} = A^\varepsilon$ for all formulas $A$.

EXAMPLE 27. *Let $A = \forall y\, (P^\perp(x) \vee P(y))$ and $b(x) = \varepsilon_y(P(x) \wedge P^\perp(y))$, then*

$$A^\varepsilon = P^\perp(x) \vee P(\varepsilon_y(P^\perp(x) \vee P(y))^{\perp \varepsilon}) = P^\perp(x) \vee P(b(x)).$$

*Let $a = \varepsilon_x A^\varepsilon = \varepsilon_x(P^\perp(x) \vee P(b(x)))$, then*

$$(\exists x\, A)^\varepsilon = A^\varepsilon[x \backslash a] = P^\perp(a) \vee P(b(a)).$$

In what follows we will replace an $n$-ary function symbol with an $\varepsilon$-term with $n$ free variables. In analogy to substitutions, such replacements will be written with angle brackets and $\lambda$-abstraction as $\langle f_1 \backslash \lambda \overline{x_1}.\varepsilon_{y_1} A_1, \dots, f_n \backslash \lambda \overline{x_n}.\varepsilon_{y_n} A_n \rangle$ Their application is defined by

$$
f(t_1, \dots, t_n)\sigma = \begin{cases} f(t_1\sigma, \dots, t_n\sigma), & \text{if } f \notin \mathrm{dom}(\sigma), \\ \varepsilon_y A(t_1\sigma, \dots, t_k\sigma), & \text{if } \langle f \backslash \lambda x_1 \cdots x_k.\varepsilon_y A(x_1, \dots, x_k, y)\rangle \in \sigma. \end{cases}
$$

DEFINITION 28. *Let $A$ be a sentence, let $((), s)$ be a Skolem mapping for $A$, let* q *be the position of a universal quantifier in $A$, let $\forall y\, B$ be the subformula of $A$ at* q*, and let $\overline{x} = x_1, \dots, x_n$ be the variables of the existential quantifiers dominating* q *from the outside-in. Define $\sigma_q = \langle s_q \backslash \lambda \overline{x}.\varepsilon_y B^{\perp^\varepsilon}\rangle$. If $q_1, \dots, q_m$ are the positions of universal quantifiers of $A$, define $\sigma_A = \bigcup_{i=1}^m \sigma_{q_i}$.*

In a nutshell, $\sigma_A$ replaces the Skolem symbols of $((), s)$ for $A$ (which have arbitrary names and an arity depending on the number of existential quantifiers on the path to the universal quantifier) by $\varepsilon$-terms (that include the information for which formula they provide a witness and have an arity depending on the number of free variables in that formula).

DEFINITION 29. *Let $E$ be an expansion tree. The set of critical axioms $E^\varepsilon$ of $E$ is defined inductively as*

$$
E^\varepsilon = \emptyset \text{ for an atom } E,
$$
$$
(E_1 \circ E_2)^\varepsilon = E_1^\varepsilon \cup E_2^\varepsilon \text{ for } \circ \in \{\wedge, \vee\},
$$
$$
(\exists x\, A +^{t_1} E_1 \cdots +^{t_n} E_n)^\varepsilon = \{A^\varepsilon[x \backslash t_i] \supset A^\varepsilon[x \backslash \varepsilon_x A^\varepsilon] \mid 1 \le i \le n\} \cup E_1^\varepsilon \cup \cdots \cup E_n^\varepsilon,
$$
$$
(\forall x\, A +^{s_q(\overline{t})} E_0)^\varepsilon = E_0^\varepsilon.
$$

The following lemma shows how reasoning about quantifiers is axiomatised by critical formulas over propositional logic. It is a variant of the lower part of Gentzen's mid-sequent theorem for the epsilon calculus and expansion trees for non-prenex formulas.

LEMMA 30. *Let $E$ be an expansion tree and $A = \mathrm{Sh}(E)$, then $E^\varepsilon \sigma_A \wedge \mathrm{Dp}(E)\sigma_A \supset A^\varepsilon$ is a tautology.*

*Proof.* We will show by induction that, for every subtree $E_0$ of $E$, the formula

$$
E_0^\varepsilon \sigma_A \wedge \mathrm{Dp}(E_0)\sigma_A \supset \mathrm{Sh}(E_0)^\varepsilon \sigma_A
$$

is a tautology. This suffices since $A^\varepsilon \sigma_A = A^\varepsilon$.

If $E_0$ is an atom, then $E_0^\varepsilon = \emptyset$ and $\mathrm{Dp}(E_0) = \mathrm{Sh}(E_0) = \mathrm{Sh}(E_0)^\varepsilon$. If $E_0 = E_1 \circ E_2$ for $\circ \in \{\wedge, \vee\}$, then, by induction hypothesis, both $E_1^\varepsilon \sigma_A \wedge \mathrm{Dp}(E_1)\sigma_A \supset \mathrm{Sh}(E_1)^\varepsilon \sigma_A$ and $E_2^\varepsilon \sigma_A \wedge \mathrm{Dp}(E_2)\sigma_A \supset \mathrm{Sh}(E_2)^\varepsilon \sigma_A$ are tautologies. Therefore also

$$
E_0^\varepsilon \sigma_A \wedge (\mathrm{Dp}(E_1)\sigma_A \circ \mathrm{Dp}(E_2)\sigma_A) \supset \mathrm{Sh}(E_1)^\varepsilon \sigma_A \circ \mathrm{Sh}(E_2)^\varepsilon \sigma_A
$$

is a tautology.

If $E_0$ is $\exists x\, A_0 +^{t_1} E_1 \cdots +^{t_n} E_n$, then, by the induction hypothesis, $E_i^\varepsilon \sigma_A \wedge \mathrm{Dp}(E_i)\sigma_A \supset \mathrm{Sh}(E_i)^\varepsilon \sigma_A$ is a tautology for $i = 1, \dots, n$. So, since $\mathrm{Dp}(E_0) = \bigvee_{i=1}^n \mathrm{Dp}(E_i)$, also

$$E_0^\varepsilon \sigma_A \wedge \mathrm{Dp}(E_0)\sigma_A \supset \bigvee_{i=1}^{n} \mathrm{Sh}(E_i)^\varepsilon \sigma_A$$

is a tautology. But now $\mathrm{Sh}(E_i)^\varepsilon \sigma_A = (A_0[x\backslash t_i])^\varepsilon \sigma_A = A_0^\varepsilon[x\backslash t_i]\sigma_A$. Moreover, $A_0^\varepsilon[x\backslash t_i] \supset A_0^\varepsilon[x\backslash\varepsilon_x A_0^\varepsilon] \in E_0^\varepsilon$ and $A_0^\varepsilon[x\backslash\varepsilon_x A_0^\varepsilon] = \mathrm{Sh}(E_0)^\varepsilon$, so $E_0^\varepsilon \sigma_A \wedge \mathrm{Dp}(E_0)\sigma_A \supset \mathrm{Sh}(E_0)^\varepsilon \sigma_A$ is a tautology.

If $E_0$ is $\forall x\, A_0 +^{s_q(\bar{t})} E_1$, then, by the induction hypothesis, $E_1^\varepsilon \sigma_A \wedge \mathrm{Dp}(E_1)\sigma_A \supset \mathrm{Sh}(E_1)^\varepsilon \sigma_A$ is a tautology. But now $E_0^\varepsilon = E_1^\varepsilon$, $\mathrm{Dp}(E_0) = \mathrm{Dp}(E_1)$ and since $\sigma_A$ replaces $s_q(\bar{t})$ by $\varepsilon_x A_0^\perp[\bar{y}\backslash\bar{t}]$ we also have $\mathrm{Sh}(E_0)^\varepsilon \sigma_A = (\forall x\, A_0)^\varepsilon \sigma_A = A_0[x\backslash s_q(\bar{t})]\sigma_A = \mathrm{Sh}(E_1)^\varepsilon \sigma_A$. □

EXAMPLE 31. *Let $D$ be as in Example* 11 *and let $A$ and $b(x)$ be as in Example* 27, *then*

$$\sigma_{\mathrm{Sh}(D)} = \langle f\backslash\lambda x.b(x)\rangle = \langle f\backslash\lambda x.\varepsilon_y(P(x) \wedge P^\perp(y))\rangle,$$
$$D^\varepsilon \sigma_{\mathrm{Sh}(D)} = \{A^\varepsilon[x\backslash c] \supset A^\varepsilon[x\backslash\varepsilon_x A^\varepsilon], A^\varepsilon[x\backslash b(c)] \supset A^\varepsilon[x\backslash\varepsilon_x A^\varepsilon]\},$$
$$\mathrm{Dp}(D)\sigma_{\mathrm{Sh}(D)} = P^\perp(c) \vee P(b(c)) \vee P^\perp(b(c)) \vee P(b(b(c))),$$

*and since* $A^\varepsilon[x\backslash c] = P^\perp(c) \vee P(b(c))$, $A^\varepsilon[x\backslash b(c)] = P^\perp(b(c)) \vee P(b(b(c)))$, *and* $\mathrm{Sh}(D)^\varepsilon = A^\varepsilon[x\backslash\varepsilon_x A^\varepsilon]$ *also*

$$D^\varepsilon \sigma_{\mathrm{Sh}(D)} \wedge \mathrm{Dp}(D)\sigma_{\mathrm{Sh}(D)} \supset \mathrm{Sh}(D)^\varepsilon$$

*is a tautology.*

DEFINITION 32. *Let $P_1 = (C_1; S_1, E)$ and $P_2 = (C_2; S_2, E')$ be $\varepsilon$-preproofs s.t. $S_1$ and $S_2$ are mergeable and $A = \mathrm{Sh}(E) = \mathrm{Sh}(E')^\perp$. Then we define the $\varepsilon$-preproof*

$$\mathrm{cut}_A(P_1, P_2) = (C_1\sigma_A \cup C_2\sigma_{A^\perp} \cup E^\varepsilon \sigma_A \cup E'^\varepsilon \sigma_{A^\perp}; S_1\sigma_A \sqcup S_2\sigma_{A^\perp}).$$

LEMMA 33. *Let $P_1 = (C_1; S_1, E)$ and $P_2 = (C_2; S_2, E')$ be $\varepsilon$-proofs s.t. $S_1$ and $S_2$ are mergeable and $A = \mathrm{Sh}(E) = \mathrm{Sh}(E')^\perp$. Then $\mathrm{cut}_A(P_1, P_2)$ is an $\varepsilon$-proof.*

*Proof.* Since $P_1$ and $P_2$ are $\varepsilon$-proofs, the formulas

$$\bigwedge C_1\sigma_A \supset \mathrm{Dp}(S_1)\sigma_A \vee \mathrm{Dp}(E)\sigma_A \text{ and } \bigwedge C_2\sigma_{A^\perp} \supset \mathrm{Dp}(S_2)\sigma_{A^\perp} \vee \mathrm{Dp}(E')\sigma_{A^\perp}$$

are tautologies. Furthermore, by Lemma 30, so are

$$E^\varepsilon \sigma_A \wedge \mathrm{Dp}(E)\sigma_A \supset A^\varepsilon \text{ and } E'^\varepsilon \sigma_{A^\perp} \wedge \mathrm{Dp}(E')\sigma_{A^\perp} \supset A^{\perp\varepsilon}.$$

Therefore

$$\bigwedge C_1\sigma_A \bigwedge C_2\sigma_{A^\perp} \bigwedge E^\varepsilon \sigma_A \bigwedge E'^\varepsilon \sigma_{A^\perp} \supset (\mathrm{Dp}(S_1)\sigma_A \vee A^\varepsilon) \wedge (\mathrm{Dp}(S_2)\sigma_{A^\perp} \vee A^{\perp\varepsilon})$$

is a tautology. Since $A^{\perp\varepsilon} = A^{\varepsilon\perp}$ also

$$\bigwedge C_1\sigma_A \bigwedge C_2\sigma_{A^\perp} \bigwedge E^\varepsilon \sigma_A \bigwedge E'^\varepsilon \sigma_{A^\perp} \supset (\mathrm{Dp}(S_1)\sigma_A \vee \mathrm{Dp}(S_2)\sigma_{A^\perp})$$

is a tautology. So, by Lemma 13, $\mathrm{cut}_A(P_1, P_2)$ is an $\varepsilon$-proof. □

This cut-operation on expansion proofs, together with the algorithm mentioned in Section 3 and described in detail in [7] gives an algorithm for translating a sequent calculus proof with cut to an $\varepsilon$-proof.

**§6. Associativity of cut.** It is well known that, in the sequent calculus, two cuts commute as in

$$
\cfrac{\cfrac{(P_1) \qquad (P_2)}{\cfrac{\Gamma, A_1 \quad A_1{}^\perp, \Delta, A_2}{\Gamma, \Delta, A_2} \qquad (P_3)}{\Gamma, \Delta, \Pi}}{}
\qquad \longleftrightarrow \qquad
\cfrac{(P_1)\atop\Gamma, A_1 \quad \cfrac{\cfrac{(P_2) \qquad (P_3)}{A_1{}^\perp, \Delta, A_2 \quad A_2{}^\perp, \Pi}}{A_1{}^\perp, \Delta, \Pi}}{\Gamma, \Delta, \Pi}
$$

for formulas $A_1, A_2$ and sequents $\Gamma, \Delta, \Pi$ with $A_1, A_2 \notin \Gamma$, $A_1{}^\perp, A_2 \notin \Delta$ and $A_1{}^\perp, A_2{}^\perp \notin \Pi$. From a category-theoretic point of view it is desirable to equip proofs with an associative cut operation. In this section we prove that our cut operation is indeed associative, i.e., the two above sequent calculus proofs yield the same expansion proof.

In order to show this result we need to prove a few simple commutation properties first.

LEMMA 34. *Let $A$ be a sentence. Then*:

1. *for every first-order term* $t$: $(t\sigma_A)^\varepsilon = t^\varepsilon \sigma_A$,
2. *for every formula $F$*: $(F\sigma_A)^\varepsilon = F^\varepsilon \sigma_A$, *and*
3. *for every expansion tree $E$ s.t. the ranges of the Skolem mappings of $E$ and $A$ are disjoint*: $(E\sigma_A)^\varepsilon = E^\varepsilon \sigma_A$.

*Proof.* 1 can be shown by induction on the structure of $t$. 2 can be shown by induction on the structure of $F$ using 1 for atoms. 3 is shown by induction on the structure of $E$: the cases of atoms, conjunction, and disjunction are straightforward. If $E$ is $\forall x\, B +^{s_q(\bar{t})} E_0$ then

$$
(E\sigma_A)^\varepsilon = (\forall x\, B\sigma_A +^{s_q(\bar{t}\sigma_A)} E_0\sigma_A)^\varepsilon = (E_0\sigma_A)^\varepsilon = E_0^\varepsilon \sigma_A = E^\varepsilon \sigma_A.
$$

If $E$ is $\exists x\, B +^{t_1} E_1 \cdots +^{t_n} E_n$ then

$$
(E\sigma_A)^\varepsilon = \{(B\sigma_A)^\varepsilon[x\backslash t_i\sigma_A] \supset (B\sigma_A)^\varepsilon[x\backslash \varepsilon_x(B\sigma_A)^\varepsilon] \mid 1 \le i \le n\} \cup (E_1\sigma_A)^\varepsilon \cup \cdots \cup (E_n\sigma_A)^\varepsilon
$$

and

$$
E^\varepsilon \sigma_A = \{B^\varepsilon[x\backslash t_i]\sigma_A \supset B^\varepsilon[x\backslash \varepsilon_x B^\varepsilon]\sigma_A \mid 1 \le i \le n\} \cup E_1^\varepsilon \sigma_A \cup \cdots E_n^\varepsilon \sigma_A.
$$

By 2 we have $(B\sigma_A)^\varepsilon[x\backslash t_i\sigma_A] = B^\varepsilon \sigma_A[x\backslash t_i\sigma_A] = B^\varepsilon[x\backslash t_i]\sigma_A$ and $(B\sigma_A)^\varepsilon[x\backslash \varepsilon_x (B\sigma_A)^\varepsilon] = B^\varepsilon \sigma_A[x\backslash \varepsilon_x B^\varepsilon \sigma_A] = B^\varepsilon[x\backslash \varepsilon_x B^\varepsilon]\sigma_A$ which entails the claim together with the induction hypothesis. $\square$

THEOREM 35. *Let $P_1 = (C_1; S_1, E_1)$, $P_2 = (C_2; E_1', S_2, E_2)$, and $P_3 = (C_3; E_2', S_3)$ be $\varepsilon$-preproofs, let $A_1 = \mathrm{Sh}(E_1) = \mathrm{Sh}(E_1')^\perp$, and let $A_2 = \mathrm{Sh}(E_2) = \mathrm{Sh}(E_2')^\perp$ s.t. $A_1, A_2 \notin \mathrm{Sh}(S_1)$, $A_1{}^\perp, A_2 \notin \mathrm{Sh}(S_2)$ and $A_1{}^\perp, A_2{}^\perp \notin \mathrm{Sh}(S_3)$. Then*

$$
\mathrm{cut}_{A_2}(\mathrm{cut}_{A_1}(P_1, P_2), P_3) = \mathrm{cut}_{A_1}(P_1, \mathrm{cut}_{A_2}(P_2, P_3)).
$$

*Proof.* Note that $A_i\sigma_{A_j} = A_i\sigma_{A_j^\perp} = A_i$ for $i, j \in \{1, 2\}$. Therefore, by definition, we have

$$
\mathrm{cut}_{A_2}(\mathrm{cut}_{A_1}(P_1, P_2), P_3) = (C_1\sigma_{A_1}\sigma_{A_2\sigma_{A_1^\perp}} \cup C_2\sigma_{A_1^\perp}\sigma_{A_2\sigma_{A_1^\perp}} \cup C_3\sigma_{A_2^\perp} \cup
$$
$$
E_1^\varepsilon \sigma_{A_1}\sigma_{A_2\sigma_{A_1^\perp}} \cup E_1'^\varepsilon \sigma_{A_1^\perp}\sigma_{A_2\sigma_{A_1^\perp}} \cup (E_2\sigma_{A_1^\perp})^\varepsilon \sigma_{A_2\sigma_{A_1^\perp}} \cup E_2'^\varepsilon \sigma_{A_2^\perp};
$$
$$
S_1\sigma_{A_1}\sigma_{A_2\sigma_{A_1^\perp}} \sqcup S_2\sigma_{A_1^\perp}\sigma_{A_2\sigma_{A_1^\perp}} \sqcup S_3\sigma_{A_2^\perp})
$$

$$\mathrm{cut}_{A_1}(P_1, \mathrm{cut}_{A_2}(P_2, P_3)) = (C_1\sigma_{A_1} \cup C_2\sigma_{A_2}\sigma_{A_1^\perp\sigma_{A_2}} \cup C_3\sigma_{A_2^\perp}\sigma_{A_1^\perp\sigma_{A_2}} \cup$$
$$E_1^\varepsilon\sigma_{A_1} \cup (E_1'\sigma_{A_2})^\varepsilon\sigma_{A_1^\perp\sigma_{A_2}} \cup E_2^\varepsilon\sigma_{A_2}\sigma_{A_1^\perp\sigma_{A_2}} \cup E_2'^\varepsilon\sigma_{A_2^\perp}\sigma_{A_1^\perp\sigma_{A_2}};$$
$$S_1\sigma_{A_1} \sqcup S_2\sigma_{A_2}\sigma_{A_1^\perp\sigma_{A_2}} \sqcup S_3\sigma_{A_2^\perp}\sigma_{A_1^\perp\sigma_{A_2}})$$

Since $C_1\sigma_{A_1}$ does not contain Skolem symbols from $A_2$, we have i) $C_1\sigma_{A_1}\sigma_{A_2} = C_1\sigma_{A_1}$ and, symmetrically ii) $C_3\sigma_{A_2^\perp}\sigma_{A_1^\perp} = C_3\sigma_{A_2^\perp}$. Analogously we obtain iii) $S_1\sigma_{A_1}\sigma_{A_2} = S_1\sigma_{A_1}$, iv) $S_3\sigma_{A_2^\perp}\sigma_{A_1^\perp} = S_3\sigma_{A_2^\perp}$, v) $E_1^\varepsilon\sigma_{A_1}\sigma_{A_2} = E_1^\varepsilon\sigma_{A_1}$, and vi) $E_2'^\varepsilon\sigma_{A_2^\perp}\sigma_{A_1^\perp} = E'^\varepsilon\sigma_{A_2^\perp}$. Since $A_1$ and $A_2$ have different Skolem symbols we have vii) $\sigma_{A_1}\sigma_{A_2} = \sigma_{A_2}\sigma_{A_1}$ and analogously viii) $\sigma_{A_1^\perp}\sigma_{A_2} = \sigma_{A_2}\sigma_{A_1^\perp}$. The equality follows straightforwardly from i)–viii) and Lemma 3. □

In fact, the $\varepsilon$-proof read off from a sequent calculus proof with cut is global in the sense that it is invariant under a wide class of rule permutations. The associativity of cut is merely a special case.

**§7. A remark on Skolem axioms.** There is an intimate relationship between the epsilon calculus and Skolemisation. In fact, by considering $\varepsilon$-terms as Skolem terms, a critical formula $A[y\backslash t] \supset A[y\backslash \varepsilon_y A]$ can be considered an instance of the Skolem axiom $\forall\overline{x}(\exists y A_0 \supset A_0[y\backslash f(\overline{x})])$. Then the elimination of critical formulas corresponds to the elimination of Skolem axioms from a proof. This relationship has been made explicit in various forms in the literature: A proof of the second $\varepsilon$-theorem based on sequent calculus and Skolem functions is given in [20]. A reformulation of first-order predicate logic using a version of the $\varepsilon$-calculus that uses Skolem functions instead of $\varepsilon$'s is given in [13]. A presentation of the $\varepsilon$ substitution method for first-order arithmetic and pure first-order logic in terms of Skolem function symbols is given in [25]. In [8, 9] the authors give a reformulation of the extended first epsilon theorem using Skolem functions instead of epsilons. In [23] Mints relates a contribution of Skolem to the epsilon calculus. This connection to Skolemisation carries over to the setting of this paper: in a formalism obtained from replacing critical formulas by (instances of) Skolem axioms, an analogue of Theorem 22 can be shown with, essentially, the same proof.

**§8. Conclusion.** We have presented a simplified proof of the epsilon theorems. The use of expansion proofs eliminates the need for Skolemisation and prenexification and their respective inverse transformations in the proof of the second epsilon theorem. This work also shows that expansion proofs [22] integrate seamlessly with epsilon terms [17] and that the notion of rank applies to lambda terms, not just to $\varepsilon$-terms. A more modern presentation including the use of expansion trees, lambda tree syntax, reduction relations, negation normal forms, a more modular treatment of the substitution of $\varepsilon$-terms and a more liberal reduction leads to a smoother formalism than the original of [17].

Epsilon proofs as introduced in this paper use expansion proofs for representing the instances of the sequent being proved as do the formalisms in [4, 16] and, up to minor modifications, also [14, 21]. These formalisms contain cuts with cut formulas and represent their instances by expansion trees as well. Epsilon proofs do not contain

cuts. Their critical axioms are more general as the translation in Section 5 shows. In how far this higher generality affects proof complexity is an open question that is closely related to the complexity of deskolemisation, see Problem 22 in [12] as well as [5, 7] for partial results.

In addition, we have observed that the natural cut operation on $\varepsilon$-proofs is associative. This suggests the perspective of an investigation of the $\varepsilon$-calculus from the point of view of category theory along the lines of [18].

## BIBLIOGRAPHY

[1] Ackermann, W. (1940). Zur Widerspruchsfreiheit der Zahlentheorie. *Mathematische Annalen*, **117**, 162–194.

[2] Aguilera, J. P., & Baaz, M. (2019). Unsound inferences make proofs shorter. *Journal of Symbolic Logic*, **84**(1), 102–122.

[3] Andrews, P. B., Bishop, M., Issar, S., Nesmith, D., Pfenning, F., & Xi, H. (1996). TPS: A theorem-proving system for classical type theory. *Journal of Automated Reasoning*, **16**(3), 321–353.

[4] Aschieri, F., Hetzl, S., & Weller, D. (2019). Expansion trees with cut. *Mathematical Structures in Computer Science*, **29**(8), 1009–1029.

[5] Avigad, J. (2003). Eliminating definitions and Skolem functions in first-order logic. *ACM Transactions on Computational Logic*, **4**, 402–415.

[6] Avigad, J., & Zach, R. The epsilon calculus. In Zalta, E. N., editor. *The Stanford Encyclopedia of Philosophy* (Spring 2020 Edition). Stanford: Metaphysics Research Lab, Stanford University, 2020.

[7] Baaz, M., Hetzl, S., & Weller, D. (2012). On the complexity of proof deskolemization. *Journal of Symbolic Logic*, **77**(2), 669–686.

[8] Baaz, M., Leitsch, A., & Lolic, A. (2018). A sequent-calculus based formulation of the extended first epsilon theorem. In Artëmov, S. N. and Nerode, A., editors. *International Symposium on the Logical Foundations of Computer Science (LFCS)*. Lecture Notes in Computer Science, Vol. 10703. Cham: Springer, pp. 55–71.

[9] ———. (2020). An abstract form of the first epsilon theorem. *Journal of Logic and Computation*, **30**(8), 1447–1468.

[10] Chaudhuri, K., Hetzl, S., & Miller, D. (2016). A multi-focused proof system isomorphic to expansion proofs. *Journal of Logic and Computation*, **26**(2), 577–603.

[11] Church, A. (1940). A formulation of the simple theory of types. *Journal of Symbolic Logic*, **5**(2), 56–68.

[12] Clote, P., & Krajíček, J. (1993). Open problems. In Clote, P. and Krajíček, J., editors. *Arithmetic, Proof Theory and Computational Complexity*. Oxford: Oxford University Press, pp. 1–19.

[13] Davis, M., & Fechter, R. (1991). A free variable version of the first-order predicate calculus. *Journal of Logic and Computation*, **1**(4), 431–451.

[14] Heijltjes, W. (2010). Classical proof forestry. *Annals of Pure and Applied Logic*, **161**(11), 1346–1366.

[15] Hetzl, S., Libal, T., Riener, M., & Rukhaia, M. Understanding resolution proofs through Herbrand's theorem. In Galmiche, D. and Dominique Larchey-

Wendling, editors. *Automated Reasoning with Analytic Tableaux and Related Methods*. Lecture Notes in Computer Science, Vol. 8123. Berlin–Heidelberg: Springer, 2013, pp. 157–171.

[16] Hetzl, S., & Weller, D. (2013). Expansion trees with cut. Preprint. Available from: http://arxiv.org/abs/1308.0428.

[17] Hilbert, D., & Bernays, P. (1939). *Grundlagen der Mathematik II*. Berlin: Springer.

[18] Hyland, J. M. E. (2002). Proof theory in the abstract. *Annals of Pure and Applied Logic*, **114**, 43–78.

[19] Leitsch, A., & Lolic, A. (2019). Extraction of expansion trees. *Journal of Automated Reasoning*, **62**(3), 393–430.

[20] Maehara, S. (1955). The predicate calculus with $\varepsilon$-symbol. *Journal of the Mathematical Society of Japan*, **7**(4): 323–344.

[21] McKinley, R. (2013). Proof nets for Herbrand's theorem. *ACM Transactions on Computational Logic*, **14**(1), 5:1–5:31.

[22] Miller, D. (1987). A compact representation of proofs. *Studia Logica*, **46**(4), 347–370.

[23] Mints, G. (1996). Thoralf Skolem and the epsilon substitution method for predicate logic. *Nordic Journal of Philosophical Logic*, **1**(2), 133–146.

[24] Moser, G., & Zach, R. (2006). The epsilon calculus and Herbrand complexity. *Studia Logica*, **82**(1), 133–155.

[25] Tait, W. (2010). The substitution method revisited. In Feferman, S., and Sieg, W., editors. *Proofs, Categories and Computations: Essays in Honor of Grigori Mints.* London: College Publications, pp. 231–241.

[26] Weller, D. (2011). On the elimination of quantifier-free cuts. *Theoretical Computer Science*, **412**(49), 6843–6854.

INSTITUTE OF DISCRETE MATHEMATICS AND GEOMETRY
TU WIEN, WIEDNER HAUPTSTRASSE 8-10
1040 VIENNA, AUSTRIA
*E-mail*: stefan.hetzl@tuwien.ac.at
*URL*: https://www.dmg.tuwien.ac.at/hetzl/