# Oberwolfach
# Preprints

BAHAREH AFSHARI, STEFAN HETZL AND
GRAHAM E. LEIGH

Herbrand's Theorem as Higher Order Recursion

## Oberwolfach Preprints (OWP)

Starting in 2007, the MFO publishes a preprint series which mainly contains research results related to a longer stay in Oberwolfach. In particular, this concerns the Research in Pairs-Programme (RiP) and the Oberwolfach-Leibniz-Fellows (OWLF), but this can also include an Oberwolfach Lecture, for example.

A preprint can have a size from 1 - 200 pages, and the MFO will publish it on its website as well as by hard copy. Every RiP group or Oberwolfach-Leibniz-Fellow may receive on request 30 free hard copies (DIN A4, black and white copy) by surface mail.

Of course, the full copy right is left to the authors. The MFO only needs the right to publish it on its website *www.mfo.de* as a documentation of the research work done at the MFO, which you are accepting by sending us your file.

In case of interest, please send a **pdf file** of your preprint by email to *rip@mfo.de* or *owlf@mfo.de*, respectively. The file should be sent to the MFO within 12 months after your stay as RiP or OWLF at the MFO.

There are no requirements for the format of the preprint, except that the introduction should contain a short appreciation and that the paper size (respectively format) should be DIN A4, "letter" or "article".

On the front page of the hard copies, which contains the logo of the MFO, title and authors, we shall add a running number (20XX - XX).

We cordially invite the researchers within the RiP or OWLF programme to make use of this offer and would like to thank you in advance for your cooperation.

# Herbrand's Theorem as Higher Order Recursion

Bahareh Afshari

*University of Gothenburg*

Stefan Hetzl

*TU Wien*

Graham E. Leigh

*University of Gothenburg*

We provide a means to compute Herbrand disjunctions directly from sequent calculus proofs with cuts. Our approach associates to a first-order classical proof $\pi \vdash \exists v F$, where $F$ is quantifier free, an acyclic higher order recursion scheme $\mathscr{H}$ whose language is finite and yields a Herbrand disjunction for $\exists v F$. More generally, we show that the language of $\mathscr{H}$ contains the Herbrand disjunction implicit in any cut-free proof obtained from $\pi$ via a sequence of Gentzen-style cut reductions that always reduce the weak side of a cut before the strong side.

# Contents

# 1 Introduction

The property of being a valid first-order formula is intimately tied to the consideration of the ground, i.e., variable-free, instances of that formula. This connection is apparent in most, if not all, proofs of the completeness theorem which, in one way or another, rely on the construction of a term model. It is plainly visible in Herbrand's theorem which states that a formula is valid if, and only if, there is a finite expansion (of existential quantifiers to disjunctions and universal quantifiers to conjunctions of instances). This feature of classical first-order logic is in contrast to both classical second-order logic, whose standard semantics goes beyond the ground instances of a countable language, and intuitionistic first-order logic, which exhibits a more complicated interaction between quantifiers and propositional connectives.

Proof-theoretically, the use of instances of a formula naturally leads to analytic, cut-free, proofs. Gentzen's mid-sequent theorem makes the close connection between Herbrand expansions and cut-free proofs apparent. Taking this perspective on the cut-elimination theorem, and thereby keeping the well-known complexity bounds in mind, shows that, in essence, cut-elimination consists of the computation of a Herbrand expansion. One may ask, however, whether given a proof with cut it is possible to compute a Herbrand expansion in a more direct way, circumventing the cumbersome process of cut-elimination. There is a number of formalisms that do just that, the historically first being Hilbert's $\varepsilon$-calculus [31] (see [35] for a contemporary exposition of the $\varepsilon$-theorems in English). In [20], Gerhardy and Kohlenbach adapt Shoenfield's variant [41] of Gödel's Dialectica interpretation [21, 5] to a system of pure predicate logic. Recent work, related to proof nets, is that of Heijltjes [22] and McKinley [33], and a similar approach, in the formalism of expansion trees [34], can be found in [30]. A different method with similar aims is cut-elimination by resolution [8].

The present work is motivated by the follow-up question: what is a minimal amount of information required for computing a Herbrand expansion from a proof with cuts? An approach which has been partially successful in answering the question is the representation of proofs as tree grammars, introduced in [24] for proofs with $\Pi_1$-cuts and extended to $\Pi_2$-cuts in [1, 2]. This emphasis on minimality plays a crucial role for several applications, such as cut-introduction [26, 25, 32], inductive theorem proving [15] and the confluence behaviour of cut-elimination [28, 29, 2]. For instance, in the case of cut-introduction, which is an attempt to automatically compress proofs via introduction of cuts (i.e. lemmas), the algorithm proceeds in two steps: 1. a smallest grammar which represents a given Herbrand expansion is computed and 2. this grammar is translated to a proof with cuts. The minimality of the grammar formalism makes the first step feasible and the second step total. Algorithms for cut-introduction and inductive theorem proving are currently being implemented in the GAPT-system [16], see e.g. [27, 17]. A further theoretical application of proof grammars is in the area of proof complexity, where lower bounds on the length of proofs with cuts (which are notoriously difficult to control) are obtained by transferring lower bounds on the size of the corresponding grammar [14, 13]. Proving these lower bounds on the size of grammars is considerably simplified by them containing only a minimal amount of information.

Continuing this research effort, we demonstrate how Herbrand expansions can be represented as languages of higher order recursion schemes derived directly from first-

order proofs with cut. Higher order recursion schemes (see e.g. [37]) are a generalisation of regular tree grammars (which correspond to order-0 recursion schemes) to finite types. The representation we outline involves interpreting inference rules of proofs as non-terminals whose production rules follow the local instantiation structure of quantifiers. The type of a non-terminal is determined entirely by the quantifier complexity of the formulæ occurring in the corresponding inference, with an inference deriving $\Sigma_n \cup \Pi_n$ formulæ being represented by a non-terminal of order $n$. Cut corresponds to composition of non-terminals, and instances of contraction give rise to non-deterministic production rules. The language of the recursion scheme induces a Herbrand expansion for the end-sequent of the proof. At the level of $\Pi_2$-cuts, the schemes closely resemble the grammars introduced in [4]. The generic case of the representation, which permits capturing cuts of arbitrary quantifier complexity, turns out to be at the level of $\Pi_3$ where both sides of a cut feature $\exists\forall$ quantifier alternations.

As far as the authors are aware, the present work marks the first method of Herbrand extraction that operates directly on sequent calculus proofs. The main result can be summarised as follows and was announced in [3].

**Theorem 1.1.** *Let $F$ be a quantifier-free formula and $\pi$ a first-order proof of $\exists\vec{v}F$ in which cut-formulæ are prenex $\Pi_n$ or $\Sigma_n$. There exists an acyclic order $n$ recursion scheme $\mathscr{H}$ with language $L(\mathscr{H})$ such that: i) $\bigvee_{\vec{t}\in L(\mathscr{H})} F(\vec{t})$ is valid; ii) $|L(\mathscr{H})| \leq 2_{n+2}^{4|\pi|^3}$ where $|\pi|$ is the number of inference rules in $\pi$; iii) $L(\mathscr{H})$ contains the Herbrand set extracted from any cut-free proof that can be obtained from $\pi$ via a sequence of Gentzen-style cut reductions that always reduces to the weak (quantifier) side of a cut before the strong side.*

## 2 Sequent Calculus for Classical First-order Logic

Terms and formulæ of first-order logic are defined as usual using the connectives $\wedge$, $\vee$ and quantifiers $\forall$, $\exists$, as well as a selection of predicate and function symbols. We assume two sets of variable symbols, *free* variables, denoted $\alpha$, $\beta$, etc., and *bound* variables, $v$, $w$, etc. Upper-case Roman letters, $A$, $B$, etc. denote formulæ and upper-case Greek letters $\Gamma$, $\Delta$, etc. range over *sequents*, namely finite sequences of formulæ. We abbreviate by $\Gamma, \Delta$ the concatenation of $\Gamma$ and $\Delta$; and $\Gamma, A$ is shorthand for $\Gamma, \{A\}$. The *length* of a

$$
\begin{array}{ll}
\text{Axioms:} & A, \bar{A} \qquad \text{for } A \text{ quantifier-free} \\[2ex]
\text{Inference rules:} & \vee \dfrac{\Gamma, A, B}{\Gamma, A \vee B} \qquad \wedge \dfrac{\Gamma, A \quad \Delta, B}{\Gamma, \Delta, A \wedge B} \\[3ex]
& \forall_{\vec{\alpha}} \dfrac{\Gamma, A(\vec{\alpha}/\vec{v})}{\Gamma, \forall\vec{v}A} \qquad \exists_{\vec{r}} \dfrac{\Gamma, A(\vec{r}/\vec{v})}{\Gamma, \exists\vec{v}A} \qquad \text{cut} \dfrac{\Gamma, A \quad \Delta, \bar{A}}{\Gamma, \Delta} \\[3ex]
& \text{w} \dfrac{\Gamma}{\Gamma, A} \qquad \text{c} \dfrac{\Gamma, A, A}{\Gamma, A} \qquad \text{p} \dfrac{\Gamma, B, A, \Delta}{\Gamma, A, B, \Delta}
\end{array}
$$

Figure 1: Axioms and rules of sequent calculus

sequent $\Gamma$ is denoted $|\Gamma|$. As the order of the formulæ in a sequent is often (though not always) unimportant, we will frequently identify sequents with (finite) multisets. We write $\bar{A}$ to denote the dual of the formula $A$ obtained by de Morgan laws. Given a sequence of variable symbols $\vec{v} = (v_0, \ldots, v_{k-1})$ of length $k$, we write $\forall \vec{v} A$ and $\exists \vec{v} A$ as shorthand for $\forall v_0 \cdots \forall v_{k-1} A$ (resp. $\exists v_0 \cdots \exists v_{k-1} A$). If $\vec{t} = (t_0, \ldots, t_{k-1})$ is a sequence of terms of the same length, $A(\vec{t}/\vec{v})$ is the formula obtained from $A$ by replacing each $v_i$ by the corresponding term $t_i$, where bound variables in $A$ are renamed as necessary to avoid variable capture.

The following abbreviations will be used in later sections. For a formula $A$, we write $A_{qf}$ to indicate that $A$ is quantifier-free, and $u(A)$ (resp. $e(A)$) for the number of consecutive universal (existential) quantifiers in $A$ before encountering an existential (universal) quantifier:

$$u(\forall v A) = u(A) + 1 \qquad\qquad e(\exists v A) = e(A) + 1$$
$$u(\exists v A) = u(A_{qf}) = 0 \qquad\qquad e(\forall v A) = e(A_{qf}) = 0$$

For notational simplicity, we work in one-sided sequent calculus with explicit structural rules for weakening ($\mathsf{w}$), contraction ($\mathsf{c}$) and permutation ($\mathsf{p}$), though the results presented apply equally to two-sided (so-called Gentzen-style) sequent calculi and either form of calculus without explicit structural rules. The axioms and rules of the calculus are laid out in Figure 1. The quantifier introduction rules $\forall_{\vec{\alpha}}$ and $\exists_{\vec{r}}$ introduce a sequence of quantifiers in one application. Applications of $\forall_{\vec{\alpha}}$ are subject to an eigenvariable condition that if $\vec{\alpha} = (\alpha_0, \ldots, \alpha_{k-1})$ then $\alpha_i$ does not occur in the sequent $\Gamma, A$ for any $i < k$. In each inference rule, the formulæ which are explicitly mentioned in the premise(s) (usually the right-most formula in the sequent) are said to be *active* in the rules applied. For example, $A$ and $B$ are active in $\wedge$ rule, both copies of $A$ are active in contraction, and there are no active formulæ in the weakening rule. Active formulæ of $\mathsf{cut}$ are refereed to as *cut formulæ*. We often leave the applications of the permutation rule implicit, writing, for instance,

$$\forall_{\vec{\alpha}} \frac{\Gamma, A(\vec{\alpha}/\vec{v}), \Delta}{\Gamma, \forall \vec{v} A, \Delta} \qquad\qquad \mathsf{cut} \frac{\Gamma, A, \Gamma' \quad \Delta, \bar{A}, \Delta'}{\Gamma, \Gamma', \Delta, \Delta'}$$

to abbreviate derivations

$$\mathsf{p}^* \frac{\forall_{\vec{\alpha}} \frac{\mathsf{p}^* \frac{\Gamma, A(\vec{\alpha}/\vec{v}), \Delta}{\Gamma, \Delta, A(\vec{\alpha}/\vec{v})}}{\Gamma, \Delta, \forall \vec{v} A}}{\Gamma, \forall \vec{v} A, \Delta} \qquad\qquad \mathsf{cut} \frac{\mathsf{p}^* \frac{\Gamma, A, \Gamma'}{\Gamma, \Gamma', A} \quad \mathsf{p}^* \frac{\Delta, \bar{A}, \Delta'}{\Delta, \Delta', \bar{A}}}{\Gamma, \Gamma', \Delta, \Delta'}$$

where in each case $\mathsf{p}^*$ denotes a sequence of permutation inferences $\mathsf{p}$, a notation we also extend to the other structural rules.

A *proof* is a finite tree labelled by sequents obtained from the axioms and rules of the calculus with the restriction that cuts apply to prenex formulæ only. Without loss of generality, we assume all proofs are *regular*, by which we mean that:

1. each eigenvariable in the proof appears in exactly one $\forall_{\vec{\alpha}}$ inference in the proof and does not occur in any sequent outside the sub-proof of this inference,

2. if $A$ appears as the active formula of a quantifier inference $\forall$ ($\exists$) then $u(A) = 0$ (resp. $e(A) = 0$).

We write $\pi \vdash \Gamma$ to express that $\pi$ is a regular proof with $\Gamma$ being the sequent appearing at the root of $\pi$. $\mathrm{EV}(\pi)$ denotes the set of eigenvariables in a proof $\pi$, and for sequences $\vec{\alpha} = (\alpha_0, \ldots, \alpha_{k-1})$ and $\vec{t} = (t_0, \ldots, t_{k-1})$ of variable symbols and terms, $\pi^{(\vec{t}/\vec{\alpha})}$ is the result of replacing throughout the proof $\pi$ each occurrence of the variable symbol $\alpha_i$ by the term $t_i$.

## 2.1 Cut Reduction and Normal Forms

The standard cut reduction and cut permutation steps are given in Figures 2 and 3. For the sake of a concise presentation, the axioms and rules are stated with implicit permutation in place. We assume all the proofs drawn in Figures 2 and 3 are regular. Hence, in the case of contraction reduction where the sub-proof $\pi_1$ is duplicated it is assumed that the eigenvariables are renamed in the copy, which is emphasised by annotating the sub-proof with an asterisk i.e. $\pi_1^*$. In the two reductions of Figure 3, r represents an arbitrary unary or binary inference rule. An example of the binary inference permutation rule for r = cut is

$$
\mathrm{cut} \cfrac{ \mathrm{cut} \cfrac{ \pi_0 \quad \pi_1 \\ \Gamma, B \qquad \Delta, \bar{B}, A }{ \Gamma, \Delta, A } \qquad \pi_2 \\ \Lambda, \bar{A} }{ \Gamma, \Delta, \Lambda }
\quad \rightsquigarrow \quad
\mathrm{cut} \cfrac{ \pi_0 \\ \Gamma, B \qquad \mathrm{cut} \cfrac{ \pi_1 \quad \pi_2 \\ \Delta, \bar{B}, A \qquad \Lambda, \bar{A} }{ \Delta, \bar{B}, \Lambda } }{ \Gamma, \Delta, \Lambda }
$$

For proofs $\pi$ and $\pi'$ we write $\pi \rightsquigarrow \pi'$ to express that $\pi'$ is obtained from $\pi$ by application of a reduction or permutation rule to a sub-proof of $\pi$, and let $\rightsquigarrow^*$ denote the reflexive transitive closure of $\rightsquigarrow$. If $\pi \rightsquigarrow \pi'$ then the reduced cut either no longer exists, is replaced by cuts on formulæ with either lower logical complexity or fewer applied contractions, or permuted to a subproof. In any given proof there may, however, be many cuts and eliminating one can (through duplicating a sub-proof) result in introducing several copies of other cuts. To obtain a cut-free proof, it is necessary to provide a (terminating) *cut elimination strategy* i.e. a procedure that given any proof $\pi \vdash \Gamma$ induces a sequence of cut reduction and permutation steps $\pi \rightsquigarrow^* \pi'$ such that $\pi' \vdash \Gamma$ and the rule cut is not used in $\pi'$.

**Theorem 2.1** (Gentzen's Hauptsatz). *There is a cut elimination strategy that transforms any proof in first-order logic to a cut-free proof.*

There are many cut elimination strategies such as top-most reduction strategy or the elimination of the cut with highest logical complexity. Different strategies provide different cut-free proofs, commonly also referred to as *normal forms*. In fact, there exist proofs with infinitely many normal forms (see e.g. [43, Example 2.1.3]). We now turn to the relationship between cut-elimination and Herbrand disjunctions in first-order logic. In the remainder of this article a *quasi cut-free proof* refers to a proof in which the only cuts are on quantifier-free formulae.
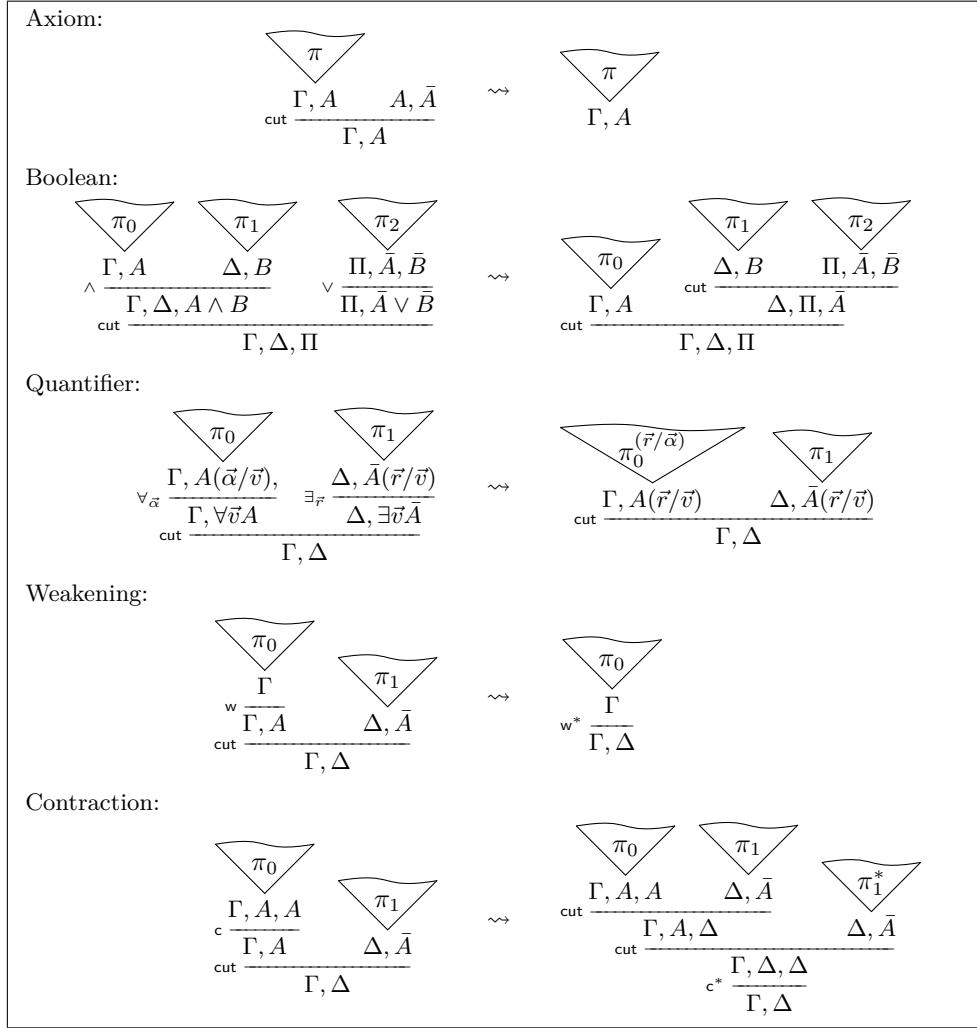
6

Axiom:

$$\text{cut}\ \dfrac{\pi \quad \Gamma, A \qquad A, \bar{A}}{\Gamma, A} \qquad \rightsquigarrow \qquad \dfrac{\pi}{\Gamma, A}$$

Boolean:

$$\text{cut}\ \dfrac{\wedge\ \dfrac{\pi_0 \quad \Gamma, A \qquad \pi_1 \quad \Delta, B}{\Gamma, \Delta, A \wedge B} \qquad \vee\ \dfrac{\pi_2 \quad \Pi, \bar{A}, \bar{B}}{\Pi, \bar{A} \vee \bar{B}}}{\Gamma, \Delta, \Pi}$$

$$\rightsquigarrow \qquad \text{cut}\ \dfrac{\pi_0 \quad \Gamma, A \qquad \text{cut}\ \dfrac{\pi_1 \quad \Delta, B \qquad \pi_2 \quad \Pi, \bar{A}, \bar{B}}{\Delta, \Pi, \bar{A}}}{\Gamma, \Delta, \Pi}$$

Quantifier:

$$\text{cut}\ \dfrac{\forall_{\vec{\alpha}}\ \dfrac{\pi_0 \quad \Gamma, A(\vec{\alpha}/\vec{v}),}{\Gamma, \forall \vec{v} A} \qquad \exists_{\vec{r}}\ \dfrac{\pi_1 \quad \Delta, \bar{A}(\vec{r}/\vec{v})}{\Delta, \exists \vec{v} \bar{A}}}{\Gamma, \Delta}$$

$$\rightsquigarrow \qquad \text{cut}\ \dfrac{\pi_0^{(\vec{r}/\vec{\alpha})} \quad \Gamma, A(\vec{r}/\vec{v}) \qquad \pi_1 \quad \Delta, \bar{A}(\vec{r}/\vec{v})}{\Gamma, \Delta}$$

Weakening:

$$\text{cut}\ \dfrac{\text{w}\ \dfrac{\pi_0 \quad \Gamma}{\Gamma, A} \qquad \pi_1 \quad \Delta, \bar{A}}{\Gamma, \Delta} \qquad \rightsquigarrow \qquad \text{w}^*\ \dfrac{\pi_0 \quad \Gamma}{\Gamma, \Delta}$$

Contraction:

$$\text{cut}\ \dfrac{\text{c}\ \dfrac{\pi_0 \quad \Gamma, A, A}{\Gamma, A} \qquad \pi_1 \quad \Delta, \bar{A}}{\Gamma, \Delta}$$

$$\rightsquigarrow \qquad \text{c}^*\ \dfrac{\text{cut}\ \dfrac{\text{cut}\ \dfrac{\pi_0 \quad \Gamma, A, A \qquad \pi_1 \quad \Delta, \bar{A}}{\Gamma, A, \Delta} \qquad \pi_1^* \quad \Delta, \bar{A}}{\Gamma, \Delta, \Delta}}{\Gamma, \Delta}$$

Figure 2: One-step cut reduction rules.

Unary inf.:

$$\text{cut}\ \dfrac{\text{r}\ \dfrac{\pi_0 \quad \Gamma', A}{\Gamma, A} \qquad \pi_1 \quad \Delta, \bar{A}}{\Gamma, \Delta}$$

$$\rightsquigarrow \qquad \text{r}\ \dfrac{\text{cut}\ \dfrac{\pi_0 \quad \Gamma', A \qquad \pi_1 \quad \Delta, \bar{A}}{\Gamma', \Delta}}{\Gamma, \Delta}$$

Binary inf.:

$$\text{cut}\ \dfrac{\text{r}\ \dfrac{\pi_0 \quad \Gamma' \qquad \pi_1 \quad \Delta', A}{\Gamma, \Delta, A} \qquad \pi_2 \quad \Lambda, \bar{A}}{\Gamma, \Delta, \Lambda}$$

$$\rightsquigarrow \qquad \text{r}\ \dfrac{\pi_0 \quad \Gamma' \qquad \text{cut}\ \dfrac{\pi_1 \quad \Delta', A \qquad \pi_2 \quad \Lambda, \bar{A}}{\Delta', \Lambda}}{\Gamma, \Delta, \Lambda}$$
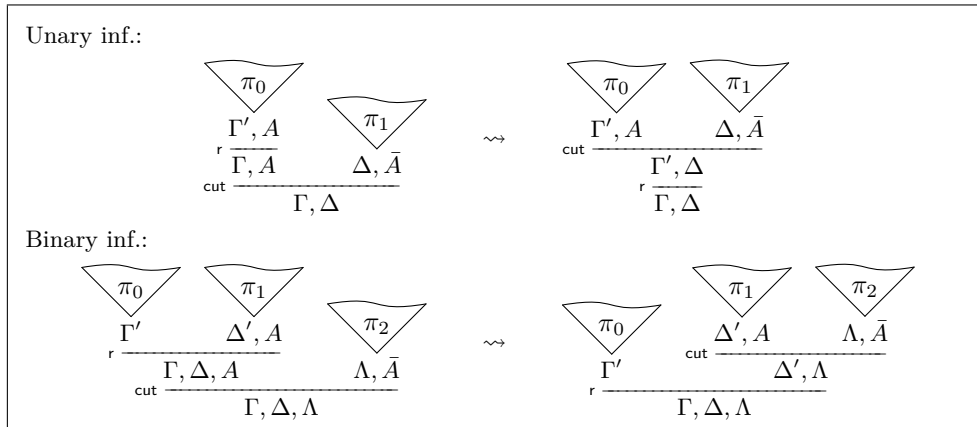
Figure 3: One-step cut permutation rules.

7

## 2.2 Herbrand's Theorem and Cut Elimination

Herbrand's theorem is considered a classic result in proof theory. It can be thought of as reducing validity in first-order logic to validity in propositional logic. From the modern perspective it can also be seen as extracting computational content to first-order proofs. A simple case of the theorem is the following.

**Theorem 2.2** (Herbrand's theorem)**.** *A formula $\exists \vec{v} A_{qf}$ is valid if and only if there exists a finite set of sequences of terms $\{\vec{t_0}, \vec{t_1}, \ldots, \vec{t_k}\}$ such that $\bigvee_{i=0}^{k} A(\vec{t_i}/\vec{v})$ is valid.*

If a formula $\exists \vec{v} A_{qf}$ is valid then any set of terms $\{\vec{t_0}, \vec{t_1}, \ldots, \vec{t_k}\}$ that validate the disjunction $\bigvee_{i=0}^{k} A(\vec{t_i}/\vec{v})$ is called a *Herbrand set*, and the disjunction itself a *Herbrand disjunction* for the formula.

Herbrand's theorem pre-dates Gentzen's Hauptsatz but the latter readily provides an instructive proof of the theorem: Suppose $\exists \vec{v} A_{qf}$ is valid and fix a quasi cut-free proof $\pi \vdash \exists \vec{v} A$. It is possible to permute the rules applied in $\pi$ so that no quantifier inference occurs above a purely propositional rule (Gentzen's mid-sequent theorem [18]). Once the proof is partitioned into a *propositional part* and a *quantifier part*, the terms that validate the formula can be directly read off from the *mid-sequent*, the sequent separating the two parts.

Herbrand's original statement is much more general than that stated above and applies to any formula of first-order logic thanks to *Herbrandisation*, the dual notion of Skolemization. Given an arbitrary formula $A$, by introducing suitable constant and function symbols it is possible to remove universal quantifiers in $A$ and obtain a $\Sigma_1$ prenex-formula which is equi-valid to $A$. Herbrandisation can also be applied to a proof of a sequent $\Gamma$ transforming it to a proof of the Herbandisation of $\Gamma$.

If a Herbrand set (disjunction) is obtained via cut elimination it is customary to refer to it as a Herbrand set (disjunction) of the proof. Note that these are not unique: different reduction strategies can lead to non-elementary many pairwise distinct Herbrand disjunctions [6]. For both computing and representing Herbrand disjunctions it is therefore desirable to bypass cut elimination. There has been a number of successful approaches such as via Herbrand nets [33], proof forests [22], expansion trees with cut [30] and functional interpretation [20]. In the next section we introduce a fresh approach using higher order recursion schemes which allows the extraction of Herbrand disjunctions directly from proofs in sequent calculus and represents them in a standard formalism from formal language theory.

## 3 Recursion Schemes for First-order Proofs

In this section we associate to each sequent calculus proof $\pi$ with $\Sigma_1$ end-sequent a non-deterministic higher order recursion scheme $\mathscr{H}_\pi$. We begin with definition of the type system and terms that will be used throughout the paper. In sections 3.2 and 3.3, higher order recursion schemes over the type system are introduced and upper bounds on the size of languages of acyclic schemes are established. The definition of $\mathscr{H}_\pi$ is given in section 3.4.

## 3.1 Types and Terms

The type system we utilise extends the hierarchy of simple types (over a type of individuals $\iota$) by pair types and two additional type constants. These are the unit type, denoted $\epsilon$, and a type $\varsigma$ of (stacks of) substitutions, elements of which are finite sequences of pairs $(\alpha, r)$ where $\alpha$ and $r$ are elements of some (and the same) type. We are interested specifically in the case that $\alpha$ is a constant symbol (of simple type) from a particular ranked alphabet $\Sigma$, and refer to the type $\varsigma$ as the type of *substitution stacks (over $\Sigma$)*, or simply $\Sigma$-*substitutions*.

The informal reading behind the type $\varsigma$ is that of an accumulator for a sequence of substitutions that are generated by reading a particular thread through a formal proof: when a witness to an existential quantifier is encountered along such a thread, the witness is outputted accompanied by the current stack of substitutions. The substitutions are not evaluated at the formal level but recorded as an element of $\varsigma$.

We begin with a formal definition of the types and conventions for their representation, followed by ranked alphabets and the recursive definition of (typed) terms including the precise form of inhabitants of the type of substitution stacks.

**Definition 3.1.** The *types* are defined in the following way.

- $\iota$ is a type, called the type of *individuals*.

- $\epsilon$ is a type, called the *unit type*.

- $\varsigma$ is a type, called the type of *substitution stacks*.

- *Function types*: if $\rho$, $\sigma$ are types then $\rho \to \sigma$ is a type.

- *Pair types*: if $\rho$, $\sigma$ are types then $\rho \times \sigma$ is a type.

A type formed without reference to $\varsigma$ is called *basic*, and one formed only out of $\iota$ and $\to$ is *simple*. The types $\iota$ and $\epsilon$ are referred to collectively as *ground types* and any type that is not a function type is called *prime*. The *sequence types* are the types of the form $\iota^n$ for any $n$, where $\iota^0 = \epsilon$ and $\iota^{n+1} = \iota \times \iota^n$. The set of all types is denoted Type.

We follow the convention that the two type forming operations $\times$ and $\to$ associate to the right, and that $\to$ binds more strongly than $\times$, so for $\rho_0, \ldots, \rho_k$ types we have

$$\rho_0 \times \rho_1 \times \cdots \times \rho_k = \rho_0 \times (\rho_1 \times \cdots \times \rho_k)$$
$$\rho_0 \to \rho_1 \to \cdots \to \rho_k = \rho_0 \to (\rho_1 \to \cdots \to \rho_k)$$
$$\rho_0 \times \cdots \times \rho_i \to \rho_{i+1} \times \cdots \times \rho_k = (\rho_0 \times \cdots \times \rho_i) \to (\rho_{i+1} \times \cdots \times \rho_k)$$

Every type $\rho$ has a unique decomposition $\rho = \rho_1 \to \rho_2 \to \cdots \to \rho_k \to co(\rho)$ where $co(\rho)$ is a prime type. Given such a decomposition of $\rho$ we refer to $co(\rho)$ as the *co-domain* of $\rho$, to $k$ as the *arity* of $\rho$, and to $\rho_i$ $(1 \leq i \leq k)$ as the *i-th domain* of $\rho$.

We now define the *order* of a type which extends the usual definition of order for the simple types. Motivated by later technicalities, it is convenient to assign order $-1$ to the type of substitutions, and order $0$ to any function type with co-domain $\epsilon$.

**Definition 3.2** (Order)**.** The *order* of a type $\rho$, $ord(\rho)$, is defined as follows.

$$ord(\iota) = ord(\epsilon) = 0 \qquad ord(\varsigma) = -1 \qquad ord(\rho \times \sigma) = \max\{ord(\rho), ord(\tau)\}$$

$$ord(\rho \to \sigma) = \begin{cases} 0, & \text{if } co(\sigma) = \epsilon, \\ \max\{ord(\rho) + 1, ord(\sigma)\}, & \text{otherwise.} \end{cases}$$

**Definition 3.3** (Ranked alphabet)**.** A (ranked) *alphabet* is a pair $\mathcal{A} = \langle S, \lambda \rangle$ where $S$ is a set, called the *carrier* of $\mathcal{A}$, and $\lambda \colon S \to \text{Type}$ is a type assignment for elements of $S$. If $\lambda(\mathsf{S})$ is a simple (basic) type for every $\mathsf{S} \in S$ we call $\mathcal{A}$ *simple* (resp. *basic*). Two ranked alphabets are *disjoint* just in case their carriers are disjoint sets.

Given an alphabet $\mathcal{A} = \langle S, \lambda \rangle$, we write $\alpha^\rho \in \mathcal{A}$ if $\alpha \in S$ and $\lambda(\alpha) = \rho$, and hence frequently identify $\mathcal{A}$ with the set $\{\alpha^{\lambda(\alpha)} \mid \alpha \in S\}$ of symbols with type annotations. For alphabets $\mathcal{A} = \langle S, \lambda \rangle$ and $\mathcal{B} = \langle S', \lambda' \rangle$, we write $\mathcal{A} \subset \mathcal{B}$ if $S \subseteq S'$ and $\lambda = \lambda' \restriction S$. In case $\mathcal{A}$ and $\mathcal{B}$ are disjoint, $\mathcal{A} \cup \mathcal{B}$ denotes the alphabet formed by the union of $\mathcal{A}$ and $\mathcal{B}$, namely $\langle S \cup S', \lambda \cup \lambda' \rangle$. The empty alphabet is denoted $\emptyset$.

**Definition 3.4** (Terms and substitutions)**.** Fix alphabets $\Sigma \subset \mathcal{A}$ where $\Sigma$ is simple. The $\mathcal{A}$-*terms over* $\Sigma$ (henceforth $\mathcal{A}$-*terms*) and the types they inhabit are defined inductively as follows, where $r : \rho$ expresses that $r$ is an $\mathcal{A}$-term of type $\rho$.

1. $\langle \rangle$ is an $\mathcal{A}$-term of type $\epsilon$.

2. If $\alpha^\rho \in \mathcal{A}$ then $\alpha$ is an $\mathcal{A}$-term of type $\rho$.

3. If $r : \rho$ and $s : \sigma$ then $\langle r, s \rangle$ is a $\mathcal{A}$-term of type $\rho \times \sigma$.

4. If $r : \sigma \to \tau$ and $s : \sigma$ then $rs$ is a $\mathcal{A}$-term of type $\tau$.

5. $\bot$ is an $\mathcal{A}$-term of type $\varsigma$.

6. If $a : \varsigma$ and $r : \rho$, and $\alpha^\rho \in \Sigma$ then $[\alpha \hookleftarrow r]a$ is an $\mathcal{A}$-term of type $\varsigma$.

7. If $r : \rho$ and $a : \varsigma$ and $\rho \neq \varsigma$ then $r \cdot a$ is an $\mathcal{A}$-term of type $\rho$.

Note that $\lambda$-abstraction is not present in the term calculus, so the existence of terms of function type $\rho$ depends on the presence of $\mathcal{A}$-symbols with type $\rho_1 \to \cdots \to \rho_k \to \rho$.

In addition to the abbreviation $r : \rho$ used above, we occasionally write $r^\rho$ to express that $r$ is an $\mathcal{A}$-term of type $\rho$. We often drop mention of $\Sigma$ and $\mathcal{A}$ if they can be inferred from the context or are not important to the given setting, in which case $\mathcal{A}$-terms are referred to simply as *terms*. Terms arising from 1, 2 and 5 are called *constants*; terms arising from cases 3 and 4 are called *pairs* and *applications* respectively; terms of type $\varsigma$ are called *substitution stacks*; and terms of the form in 7 are called *explicit substitutions* (or simply *substitutions* if there is no cause for confusion). A *basic term* is any term constructed via the rules 1 to 4 only, i.e. a $\mathcal{B}$-term over $\emptyset$ for some basic alphabet $\mathcal{B}$. A term of sequence type is called a *sequence*. Application is assumed to associate to the left, and pairing and the formation rule for substitution stacks both associate to the right.

The sub-term relation is defined as usual over the basic terms, and is extended to terms containing substitutions by defining the sub-terms of $\bot$ to be $\{\bot\}$, the sub-terms

of $a = [\alpha \hookleftarrow r]b$ to be $a$ and any sub-term of $r$ or $b$, and the sub-terms of $r = s \cdot a$ to be $r$ and any sub-term of $s$ or $a$. Thus the basic terms are precisely those terms that do not have a substitution stack as a sub-term.

Given a finite sequence of terms $(r_i : \rho_i)_{i \leq k}$, let $\langle r_0, r_1, \ldots, r_k \rangle$ be the term $r_0$ if $k = 0$ and, otherwise, the pair $\langle r_0, \langle r_1, \ldots \langle r_{k-1}, r_k \rangle \cdots \rangle \rangle$ of type $\rho_0 \times \rho_1 \times \cdots \times \rho_k$. The *order* of a term is the order of its type.

**Proposition 3.5.** *If $\Sigma \subset \Sigma'$ are simple alphabets and $\mathcal{A}$ is an alphabet extending $\Sigma'$, then every $\mathcal{A}$-term over $\Sigma$ is an $\mathcal{A}$-term over $\Sigma'$.*

In addition to the term-level explicit substitutions, there is of course the usual operation of substituting given symbols by terms of corresponding type which we refer to as *implicit substitution*. Explicit substitutions can be interpreted as implicit substitutions by reading terms $r \cdot a$ as the image of $r$ under the (implicit) substitution described by $a$, a process we call *evaluation*. The following definitions explicate these two operations. Fix an alphabet $\mathcal{A}$.

**Definition 3.6** (Implicit substitution)**.** For $\mathcal{A}$-terms $r : \rho$, $t_0 : \tau_0$, $\ldots$, $t_k : \tau_k$ and distinct symbols $\alpha_0^{\tau_0}, \ldots, \alpha_k^{\tau_k} \in \mathcal{A}$, the term $r(\vec{t}/\vec{\alpha})$ is the $\mathcal{A}$-term given by simultaneously replacing every occurrence of $\alpha_i$ (for $i \leq k$) in $r$ by $t_i$, defined recursively by:

$$\beta(\vec{t}/\vec{\alpha}) = \begin{cases} \beta, & \text{if } \beta \in \mathcal{A} \text{ and } \beta \notin \{\alpha_i \mid i \leq k\}, \\ t_i, & \text{if } \beta = \alpha_i, \end{cases}$$

$$\langle \rangle(\vec{t}/\vec{\alpha}) = \langle \rangle \qquad\qquad (rs)(\vec{t}/\vec{\alpha}) = r(\vec{t}/\vec{\alpha})s(\vec{t}/\vec{\alpha})$$

$$\bot(\vec{t}/\vec{\alpha}) = \bot \qquad\qquad \langle r, s \rangle(\vec{t}/\vec{\alpha}) = \langle r(\vec{t}/\vec{\alpha}), s(\vec{t}/\vec{\alpha}) \rangle$$

$$([\beta \hookleftarrow s]a)(\vec{t}/\vec{\alpha}) = [\beta \hookleftarrow s(\vec{t}/\vec{\alpha})](a(\vec{t}/\vec{\alpha})) \qquad\qquad (r \cdot a)(\vec{t}/\vec{\alpha}) = r(\vec{t}/\vec{\alpha}) \cdot (a(\vec{t}/\vec{\alpha}))$$

If the choice of $\vec{\alpha}$ can be inferred from context, we write $r(\vec{t})$ in place of $r(\vec{t}/\vec{\alpha})$.

**Definition 3.7** (Evaluating substitutions)**.** Given an $\mathcal{A}$-term $r$ and a substitution stack $a = [\alpha_1 \hookleftarrow s_1] \cdots [\alpha_k \hookleftarrow s_k]\bot : \varsigma$ over some simple alphabet $\Sigma \subset \mathcal{A}$, the *evaluation of $r$ relative to $a$* is the $\mathcal{A}$-term over $\Sigma$ given by

$$r^a := r(s_1/\alpha_1) \cdots (s_k/\alpha_k).$$

The *evaluation* of $r$ is the term $r^\circ$ given by recursively evaluating relative to each substitution in $r$, namely evaluation leaves basic terms unchanged, commutes with application and pairing, is defined by $\bot^\circ = \bot$ and $([\alpha \hookleftarrow r]a)^\circ = [\alpha \hookleftarrow r^\circ]a^\circ$ on substitution stacks, and by $(r \cdot a)^\circ = (r^\circ)^{a^\circ}$ for explicit substitutions.

Note that the evaluation of a substitution stack on a term is well-defined due to the typing constraints on their formation.

An alphabet generally specifies a set of symbols which are associated certain re-write rules in a recursion scheme. In this context, an explicit substitution acts as a delayed substitution which is not evaluated until no further re-writes to sub-terms are possible. For instance, over the alphabet $\{\mathsf{F}^{\iota \to \iota}, \mathsf{G}^\iota, \mathsf{e}^\iota, \circ^{\iota \to \iota \to \iota}, \alpha^\iota\}$ with associated re-write rules

$\mathsf{F}x \to (\mathsf{G} \circ x) \cdot [\alpha \hookleftarrow \mathsf{e}]\bot$ (for any instantiation of $x$) and $\mathsf{G} \to \alpha$, a derivation starting from the term $t = \mathsf{F}\alpha$ is

$$t \to (\mathsf{G} \circ \alpha) \cdot ([\alpha \hookleftarrow \mathsf{e}]\bot)$$
$$\to (\alpha \circ \alpha) \cdot ([\alpha \hookleftarrow \mathsf{e}]\bot).$$

The final term evaluates to $\mathsf{e} \circ \mathsf{e}$. Attempting to read the explicit substitution implicitly leads also to the derivation

$$t \to (\mathsf{G} \circ \alpha) \cdot ([\alpha \hookleftarrow \mathsf{e}]\bot)$$
$$= \mathsf{G} \circ \mathsf{e}$$
$$\to \alpha \circ \mathsf{e}.$$

**Lemma 3.8.** *If $r : \rho$ is a $\Sigma$-term for some simple alphabet $\Sigma$ and $\rho$ is a basic type then $r^\circ$ is a basic $\Sigma$-term of type $\rho$.*

*Proof.* All substitution stacks that may occur in a term of basic type built from an alphabet of simply-typed symbols must be within the context of an explicit substitution. As evaluation replaces every explicit substitution by an implicit one, the result is a basic term of the same type. $\square$

**Lemma 3.9.** *If $r$ and $a = [\alpha \hookleftarrow s]b : \varsigma$ are $\mathcal{A}$-terms such that $\alpha$ does not occur in $r$, then $r^a = r^b$.*

**Definition 3.10** ($\Sigma$-length)**.** Given alphabets $\Sigma \subset \mathcal{A}$ and an $\mathcal{A}$-term $s$, the $\Sigma$-*length* of $s$, written $|s|_\Sigma$, is the number of occurrences of symbols in $s$ that are not $\Sigma$-terms, formally: $|\alpha|_\Sigma = 0$ if $\alpha \in \Sigma \cup \{\langle\rangle\}$; $|\alpha|_\Sigma = 1$ if $\alpha \in \mathcal{A}$ and $\alpha \notin \Sigma$; and $|r|_\Sigma = |s|_\Sigma + |t|_\Sigma$ if $r \in \{\langle s, t\rangle, st, s \cdot t, [\alpha \hookleftarrow s]t\}$.

In particular, the $\Sigma$-length of a $\Sigma$-term is 0 and if $\Sigma$ is the empty alphabet then $\Sigma$-length of any term is the number of leaves in the tree representation of the term not labelled by $\langle\rangle$.

**Remark 3.11** (Notational conventions)**.** Symbols $\rho$, $\sigma$ and $\tau$ (also with indices) range over types. We commonly notate alphabets by upper-case Roman symbols in calligraphic typeface: $\mathcal{A}$, $\mathcal{B}$, etc, though Greek symbols $\Sigma$ and $\Sigma'$ will be used for simple alphabets. Sans-serif typeface ($\mathsf{f}$, $\mathsf{F}$, $\mathsf{s}$, $\mathsf{S}$, etc.) and lowercase Greek symbols $\alpha$, $\beta$, etc. range over elements of ranked alphabets, with the latter particularly used for constants of simple type. Italicised letters $r$, $s$, $t$, $R$, $S$, etc. range over terms and $a$, $b$ over substitution stacks, i.e. terms of type $\varsigma$.

## 3.2 Higher Order Recursion Schemes

Higher order recursion schemes provide a handy framework for extending the correspondence between formal grammars and Herbrand sets established in [1, 24]. Their advantage over formal grammars can be appreciated by the fact that they allow us to dispense with the rigidity conditions (equality constraints) that were necessary in previous approaches.

**Definition 3.12** (Higher order recursion scheme)**.** A *(non-deterministic) higher order recursion scheme*, or simply *recursion scheme*, is a tuple $\mathscr{R} = \langle \Sigma, \mathcal{N}, \mathcal{S}, \mathcal{P} \rangle$ where $\Sigma$ is a simple alphabet, $\mathcal{N}$ is a alphabet of *non-terminals* disjoint from $\Sigma$, $\mathcal{S} \subseteq \mathcal{N}$ is a designated finite set of *starting symbols* of sequence type, $\mathcal{P}$ is a set of pairs $(\mathsf{F}^\rho, t)$, called *production rules*, such that $\mathsf{F}^\rho \in \mathcal{N}$ and $t : co(\rho)$ is a $(\Sigma \cup \mathcal{N} \cup \{x_1^{\rho_1}, \ldots, x_k^{\rho_k}\})$-term over $\Sigma$ where $x_i$ is a fresh symbol not in $\mathcal{N}$ and $\rho_i$ is the $i$-th domain of $\rho$. A production rule $(\mathsf{F}^\rho, t)$ where the arity of $\rho$ is $k$ is written as

$$\mathsf{F}x_1 \cdots x_k \rightarrow_\mathscr{R} t,$$

or $\mathsf{F}\vec{x} \rightarrow_\mathscr{R} t$. Notice that by definition the term $\mathsf{F}x_1 \cdots x_k$ is of type $co(\rho)$.

A non-terminal $\mathsf{F} \in \mathcal{N}$ of $\mathscr{R}$ is *determined* if there is a unique production rule $(\mathsf{F}^\rho, t)$ in $\mathcal{P}$. By an $\mathscr{R}$*-term* we mean a $(\Sigma \cup \mathcal{N})$-term over $\Sigma$. The *order* of $\mathscr{R}$ is the supremum over orders of the types of non-terminals of $\mathscr{R}$.

Notice that we do not require that $\mathscr{R}$ contains only finitely many non-terminals, nor that the set of start symbols is non-empty. This is for technical convenience as it allows us to consider the recursion schemes of the next section as finitely generated 'sub-schemes' of a single infinite recursion scheme. Moreover, higher order recursion schemes are traditionally presented in the context of simple types, wherein start symbols are all of type $\iota$ (and indeed a single start symbol suffices) and production rules have the form $\mathsf{F}\vec{x} \rightarrow t$ with $t : \iota$. We consider the above definition to be the natural extension of recursion schemes to accommodate non-trivial prime types.

A given non-terminal may be assigned multiple production rules, leading to non-determinism. To simplify presentation of production rules in this case we adopt the convention of writing

$$\mathsf{F}\vec{x} \rightarrow_\mathscr{R} t_0 \mid \cdots \mid t_k$$

to express that $\mathscr{R}$ contains exactly the production rule $\mathsf{F}\vec{x} \rightarrow t_i$ for each $i \leq k$, i.e. $\mathsf{F}\vec{x} \rightarrow_\mathscr{R} t_i$ for each $i \leq k$ and if $\mathsf{F}\vec{x} \rightarrow_\mathscr{R} t$ then $t = t_i$ for some $i \leq k$.

**Definition 3.13** (Derivations and language)**.** Let $\mathscr{R} = \langle \Sigma, \mathcal{N}, \mathcal{S}, \mathcal{P} \rangle$ be a higher order recursion scheme. We extend the relation $\rightarrow_\mathscr{R}$ to a relation on $\mathscr{R}$-terms defined by setting $r \rightarrow_\mathscr{R} s$ if either

- $r = \mathsf{F}r_1 \cdots r_k$ for some $\mathsf{F}^\rho \in \mathcal{N}$ with arity $k$ and there exists a production rule $\mathsf{F}\vec{x} \rightarrow_\mathscr{R} t$ such that $s = t(\vec{r}/\vec{x})$;

- $r = t(r_0/x)$, $s = t(s_0/x)$ and $r_0 \rightarrow_\mathscr{R} s_0$.

A *derivation* of $s$ from $r$ is a sequence $r = r_0 \rightarrow_\mathscr{R} \cdots \rightarrow_\mathscr{R} r_k = s$, the *length* of which is $k$. We say $s$ is *derivable from $r$* in $\mathscr{R}$, in symbols $r \rightarrow_\mathscr{R}^* s$ (or $r \rightarrow^* s$ if $\mathscr{R}$ is clear from the context), if there exists a derivation of $s$ from $r$, and $s$ is *derivable in $\mathscr{R}$* if $s$ is derivable from some $\mathsf{S} \in \mathcal{S}$. The *language* of $\mathscr{R}$, written $L(\mathscr{R})$, is the set of pairs $(\mathsf{S}, t)$ such that $\mathsf{S} \in \mathcal{S}$, $t$ is a basic $\Sigma$-term and $\mathsf{S} \rightarrow_\mathscr{R}^* t$.

**Definition 3.14.** Let $\mathscr{R} = \langle \Sigma, \mathcal{N}, \mathcal{S}, \mathcal{P} \rangle$ be a higher order recursion scheme. $\mathscr{R}$ is *finite* if $\mathcal{N}$ and $\mathcal{P}$ are both finite sets, and is *acyclic* if there exists a transitive, irreflexive relation $<$ on $\mathcal{N}$ such that for every production rule $\mathsf{F}\vec{x} \rightarrow_\mathscr{R} t$ and every non-terminal $\mathsf{G}$ occurring in $t$, $\mathsf{G} < \mathsf{F}$.

**Lemma 3.15.** *A finite acyclic recursion scheme induces a finite language.*

An upper bound on the size of the language of acyclic recursion schemes can be obtained by reducing the problem to the length of reduction sequences for the simply-typed $\lambda$-calculus. Bounds on normalisation in the simply-typed $\lambda$-calculus have been given by Schwichtenberg [39] and improved to exact bounds by Beckmann [10]. In the following we use Beckmann's result to obtain concrete bounds for acyclic recursion schemes. Let $2_0^n = n$ and $2_{k+1}^n = 2^{2_k^n}$ and extend the length function of the previous section to include $\lambda$-abstractions by setting $|\lambda x s|_\Sigma = |s|_\Sigma + 1$.

**Theorem 3.16** (Beckmann [10])**.** *Let $t$ be a term in the simply-typed $\lambda$-calculus over a simple alphabet $\Sigma$. The length of any $\beta$-reduction sequence starting from $t$ is bounded by $2_{d(t)}^{|t|_\Sigma}$ where $d(t)$ denotes the maximum among orders of sub-terms of $t$.*

Beckmann's bound still applies if $t$ is an arbitrary $\lambda$-term over the calculus of $\Sigma$-terms given in Definition 3.4 subject to the restriction that $\langle\rangle$ is the only sub-term of $t$ of type $\epsilon$ (a necessary restriction due to our non-standard definition of order). Non-deterministic reductions can also be incorporated via a fresh operator $|$ and permitting $\beta$-reductions of the form $(\lambda x. t_0|\cdots|t_k)s \to_\beta t_i(s/x)$ for each $i \leq k$. In this case the length and the function $d$ is given by $|s|t|_\Sigma = \max\{|s|_\Sigma, |t|_\Sigma\}$ and $d(s|t) = \max\{d(s), d(t)\}$. Finally, we wish to allow for so-called $\eta$-long reductions, i.e., reductions $(\lambda x_0 \cdots x_k. s)t_0 \cdots t_k \to_\beta s(\vec{t}/\vec{x})$ where $s$ is not an abstraction. Provided that only $\eta$-long reductions are permitted and each counts as one step in a $\beta$-reduction sequence, Beckmann's bound holds with the analogous change to the length function: $|\lambda \vec{x} s|_\Sigma = |s|_\Sigma + 1$ if $s$ not an abstraction.

From these observations we may deduce the following result. We restrict ourselves to recursion schemes built without the substitution stacks as this will suffice for our later use.

**Theorem 3.17.** *Let $\mathscr{R} = \langle \Sigma, \mathcal{N}, \mathcal{S}, \mathcal{P} \rangle$ be a finite acyclic order $n$ recursion scheme such that every non-terminal has basic type, and for every production rule $\mathsf{F}\vec{x} \to t$ in $\mathscr{R}$, $|t|_\Sigma < k$. The length of every derivation in $\mathscr{R}$ is bounded by $2_{n+1}^{|\mathcal{N}|(k+1)}$.*

*Proof.* Let $\mathscr{R} = \langle \Sigma, \mathcal{N}, \mathcal{S}, \mathcal{P} \rangle$ be an order $n$ recursion scheme fulfilling the requirements in the statement. Without loss of generality we may assume that $\mathcal{S}$ is a singleton, that every non-terminal is associated at least one production rule, and that $\mathcal{X}$ is the alphabet of variable symbols disjoint from both $\Sigma$ and $\mathcal{N}$ such that every term occurring in a production rule in $\mathscr{R}$ is an $(\Sigma \cup \mathcal{N} \cup \mathcal{X})$-term.

Fix an enumeration $\mathsf{F}_N^{\rho_N} < \cdots < \mathsf{F}_2^{\rho_2} < \mathsf{F}_1^{\rho_1}$ of the non-terminals of $\mathscr{R}$ according to a total ordering $(<)$ witnessing acyclicity of $\mathscr{R}$. We may assume $\mathcal{S} = \{\mathsf{F}_1\}$, so $\rho_1$ is prime. Let $\mathcal{Y} = \{y_1^{\rho_1}, y_2^{\rho_2}, \ldots, y_N^{\rho_N}\}$ be a set of fresh (and pairwise distinct) variable symbols of marked type. We define by recursion a sequence $s_1, \ldots, s_N$ of well-typed $\lambda$-terms all of type $\rho_1$ such that $s_i$ contains only the variables $y_{i+1}, \ldots, y_N$ free and the length of every derivation from $\mathsf{F}_1$ which only re-writes non-terminals $\mathsf{F}_j$ for $j \leq i$ is bounded by the length of the longest $\beta$-reduction sequence starting from $s_i$. Suppose $\{\mathsf{F}_i \vec{x} \to T_j : j \leq m\}$ is the set of production rules associated to $\mathsf{F}_i$ in $\mathscr{R}$. For each $j \leq m$, let $t_j = T_j(y_{i+1}, \ldots, y_N/\mathsf{F}_{i+1}, \ldots, \mathsf{F}_N)$ be the $(\Sigma \cup \mathcal{X} \cup \mathcal{Y})$-term resulting from $T_j$ by substituting the non-terminals $\mathsf{F}_{i+1}, \ldots, \mathsf{F}_N$ by variables $y_{i+1}, \ldots, y_N$ respectively. It follows that $|t_j|_\Sigma \leq |T_j|_\Sigma < k$. Finally, define $s_i = \lambda \vec{x}. t_0|\cdots|t_m$ if $i = 1$, and

$s_i = (\lambda y_i. s_{i-1})(\lambda \vec{x}. t_0 | \cdots | t_m)$ otherwise, and notice that $|s_N|_\Sigma \leq N(k+1)$ and the maximal order among sub-terms of $s_N$ is no greater than $n+1$. Every $\mathscr{R}$-derivation from $\mathsf{F}_1$ can be replicated as a sequence of one-step $\eta$-long $\beta$-reductions starting from $s_N$, the length of which, by Beckmann's bound, is no greater than $2_{n+1}^{N(k+1)}$. $\qquad\square$

As a corollary we obtain we obtain bounds on the size of languages.

**Corollary 3.18.** *Let $\mathscr{R}$ and $k$ be as in the previous theorem and suppose every non-terminal in $\mathscr{R}$ is associated at most two production rules. Then the size of $L(\mathscr{R})$ is bounded by $2_{n+2}^{|\mathcal{N}|(k+1)}$.*

*Proof.* Given a recursion scheme $\mathscr{R}$ all terms in $L(\mathscr{R})$ can be derived via the leftmost reduction strategy. By the previous theorem, the length of these derivations is bounded by $2_{n+1}^{|\mathcal{N}|(k+1)}$, leading to a bound of $2_{n+2}^{|\mathcal{N}|(k+1)}$ on the size of $L(\mathscr{R})$. $\qquad\square$

The bound given in Corollary 3.18 is optimal in the parameter $n$ as the next lemma demonstrates.

**Lemma 3.19.** *Let $\Sigma$ be the ranked alphabet $\{\mathsf{a}^\iota, \mathsf{b}^\iota, \mathsf{d}^{\iota \to \iota \to \iota \to \iota}\}$. There exists a sequence of acyclic higher order recursion schemes $\mathscr{R}_n = \langle \Sigma, \mathcal{N}_n, \mathcal{S}_n, \mathcal{P}_n \rangle$ such that*

1. *the order of $\mathscr{R}_n$ is $n$,*

2. *$|\mathcal{N}_n|, |\mathcal{P}_n| = O(n)$,*

3. *$\max\{|t|_\Sigma : \mathsf{F}\vec{x} \to_{\mathscr{R}_n} t\} = O(n)$,*

4. *$|L(\mathscr{R}_n)| \geq 2_{n+2}^1$.*

*Proof.* It suffices to translate Beckmann's lower bounds from [10] to the context of recursion schemes. Define $\tau_0 = \iota$ and $\tau_{i+1} = \tau_i \to \tau_i$ for each $i < \omega$. So $\tau_i$ has order and arity $i$ for each $i$. Fix $n > 0$. The recursion scheme $\mathscr{R}_n$ comprises a single start symbol $\mathsf{S}_n : \iota$ and a non-terminal $\mathsf{F}_i : \tau_i$ for each $i \leq n$. The production rules are

$$\mathsf{F}_0 \to \mathsf{a} \mid \mathsf{b} \qquad\qquad\qquad \mathsf{S}_n \to \mathsf{F}_n(\mathsf{F}_n\mathsf{F}_{n-1})\mathsf{F}_{n-2}\ldots\mathsf{F}_1\mathsf{F}_0$$
$$\mathsf{F}_1 x_0 \to \mathsf{d}\mathsf{F}_0 x_0 x_0 \qquad \mathsf{F}_{i+2}x_0 x_1 \cdots x_{i+1} \to x_0(x_0 x_1)x_2 \cdots x_{i+1}$$

Requirements 1–3 are clearly satisfied. To deduce 4, observe that applying deterministic production rules only, $\mathsf{S}_n \to^* \mathsf{F}_1^{(2_n^1)}\mathsf{F}_0$, where $\mathsf{X}^{(k)}$ denotes the $k$-fold iteration of $\mathsf{X}$. Thus we see that $L(\mathscr{R}_n)$ is the set of complete binary trees of height $2_n^1 + 1$ with each leaf and inner node labelled by either $\mathsf{a}$ or $\mathsf{b}$, i.e. $|L(\mathscr{R}_n)| \geq 2_{n+2}^1$. $\qquad\square$

## 3.3 Recursion Schemes with Pattern-Matching

To control the space of derivations we will utilise recursion schemes equipped with pattern-matching, introduced in [36]. In their full generality pattern-matching recursion schemes form a Turing complete model of computation, though we will only employ a subclass in which pattern-matching is restricted to decomposing sequences. The following definition presents the particular schemes we utilise.

**Definition 3.20** (Pattern-matching recursion schemes)**.** A *pattern-matching recursion scheme* is a tuple $\mathscr{R} = \langle \Sigma, \mathcal{N}, \mathcal{S}, \mathcal{P} \rangle$ where $\Sigma$, $\mathcal{N}$ and $\mathcal{S}$ are as in Definition 3.12 and $\mathcal{P}$ may include type-preserving production rules of the form

$$\mathsf{F} x_0 \cdots x_{k-1} \langle x_k, \ldots, x_{k+l} \rangle \to_{\mathscr{R}} t$$

where $t$ is a $\Sigma \cup \mathcal{N} \cup \{x_i \mid i \leq k+l\}$-term over $\Sigma$ of prime type.

The associated reduction relation $r \to_{\mathscr{R}} s$ is defined by the two conditions in Definition 3.13 and an additional clause:

- $r = \mathsf{F} r_0 \cdots r_{k-1} \langle r_k, \ldots, r_{k+l} \rangle$ for some $\mathsf{F} \in \mathcal{N}$ of arity $k$ and terms $\vec{r} = (r_i)_{i \leq k+l}$, and there exists a production rule $\mathsf{F} x_0 \cdots x_{k-1} \langle x_k, \ldots, x_{k+l} \rangle \to_{\mathscr{R}} t$ such that $s = t(\vec{r}/\vec{x})$.

The definition of a derivation and language for pattern-matching recursion schemes are analogous.

**Remark 3.21.** Pattern-matching recursion schemes can be simulated by higher order recursion schemes using constants representing projection functions for pairs in place of pattern-matching. In particular, the upper-bounds given by Theorem 3.17 and Corollary 3.18 apply to pattern-matching recursion schemes without change. There is, however, a subtle difference between the two in the presence of non-determinism and this will be exploited heavily in the next section. In the remainder of this paper *recursion scheme* refers to pattern-matching recursion schemes unless otherwise stated.

### 3.4 Herbrand Schemes

The recursion scheme associated to a proof $\pi$, which we call the *Herbrand scheme for $\pi$* and denote as $\mathscr{H}_\pi$, contains a non-terminal $\mathsf{N}^i_{\pi'}$ for each sub-proof $\pi' \vdash B_0, \ldots, B_k$ of $\pi$ and each $i \leq k$. The interpretation of such a non-terminal is a function which returns a witness (possibly containing explicit substitutions) for each weak quantifier in $B_i$ given input for each strong quantifier in the sequent. The arity of $\mathsf{N}^i_{\pi'}$ is $k + 2$, namely one greater than the length of the sequent: the first argument is a substitution stack and the $(j+1)$-th argument is the 'input' for the formula $B_j$. The type of $\mathsf{N}^i_{\pi'}$ depends only on the quantifier structure of the formulæ in the sequent. In particular, the types of $\mathsf{N}^i_{\pi'}$ and $\mathsf{N}^j_{\pi'}$ differ only in their co-domain. Furthermore, the reduction rules governing $\mathsf{N}^i_{\pi'}$ are determined by the final inference in $\pi'$ and choice of $i$, and re-write $\mathsf{N}^i_{\pi'}$ to a term containing non-terminals for the immediate sub-proofs of $\pi'$, so are independent of the particular starting proof. This property implies that the typing and re-write rules for a non-terminal $\mathsf{N}^i_{\pi'}$ are invariant across all Herbrand schemes $\mathscr{H}_\pi$ for which $\pi'$ is a sub-proof of $\pi$, whence we may consider two Herbrand schemes as comprising identical sets of non-terminals and production rules and differing only in the selection of start symbols.

We begin by introducing the types that occur most prominently in Herbrand schemes. To each prenex formula $F$ we assign two types, the *output* type, $\tau_F$, and the *input* type, $\tau_F^*$, representing the 'existential' and 'universal' structure of $F$ respectively. These types are determined by the quantifier structure of $F$ and are defined as follows. If $F$ is quantifier-free, $\tau_F = \tau_F^* = \epsilon$; otherwise,

$$\tau_{\forall v F} = \tau_F, \quad \tau_{\exists v F} = \begin{cases} \iota \times \tau_F, & \text{if } u(F) = 0, \\ \iota \times (\tau_{\bar{F}} \to \tau_F), & \text{if } u(F) > 0, \end{cases} \quad \tau_F^* = \begin{cases} \tau_{\bar{F}}, & \text{if } e(F) = 0, \\ \tau_F \to \tau_{\bar{F}}, & \text{if } e(F) > 0. \end{cases}$$

16

**Lemma 3.22.** *Let $F$ be a prenex formula and $\vec{v} = (v_i)_{i<k}$. Then*

1. *$\tau_F$ is a non-simple basic prime type.*

2. *$\tau_F = \tau_{F(\vec{r}/\vec{\alpha})}$ and $\tau_F^* = \tau_{F(\vec{r}/\vec{\alpha})}^*$.*

3. *If $e(F) > 0$ then $\tau_F^* = \tau_{\bar{F}}^* \to \tau_{\bar{F}}$.*

4. *$\tau_{\exists \vec{v} F} = \underbrace{\iota \times \cdots \times \iota}_{k} \times \tau_F^*$ and $\tau_{\forall \vec{v} F}^* = \underbrace{\iota \times \cdots \times \iota}_{k} \times \tau_F^*$.*

*Proof.* By definition and (for 4) induction. $\qquad\square$

**Example 3.23.** We compute the input and output types for prenex $\Pi_2$ and $\Sigma_2$ formulæ. Let $B = \exists \vec{w} A_{qf}$ and $C = \forall \vec{v} B$ where $\vec{w}$ and $\vec{v}$ have non-zero length $m$ and $n$ respectively.

$$
\begin{array}{llll}
\tau_B = \iota^m & \tau_{\bar{B}} = \epsilon & \tau_B^* = \iota^m \to \epsilon & \tau_{\bar{B}}^* = \iota^m \\
\tau_C = \iota^m & \tau_{\bar{C}} = \underbrace{\iota \times \cdots \times \iota}_{n} \times (\iota^m \to \epsilon) & \tau_C^* = \tau_{\bar{C}} & \tau_{\bar{C}}^* = \tau_{\bar{C}} \to \tau_C
\end{array}
$$

**Definition 3.24** (Herbrand scheme)**.** Fix a proof $\pi \vdash A_0, \ldots, A_k$ with $\Sigma_1$ end-sequent and let $\Sigma_\pi$ be the simple alphabet consisting of a constant symbol $\mathsf{c}$ of type $\iota$ and the function symbols and eigenvariables occurring in $\pi$ (typed accordingly). The *Herbrand scheme for $\pi$* is the higher order recursion scheme $\mathscr{H}_\pi = \langle \Sigma_\pi, \mathcal{N}_\pi, \mathcal{S}_\pi, \mathcal{P}_\pi \rangle$ with the following non-terminals and production rules.

1. A non-terminal $\mathsf{c}_\rho : \rho$ for each basic type $\rho \notin \{\iota, \epsilon\}$ that occurs as a sub-type of a type $\tau_B$ or $\tau_B^*$ for a formula $B$ occurring in $\pi$, with production rules

$$
\begin{array}{ll}
\mathsf{c}_\rho \to \langle \mathsf{c}_{\tau_0}, \mathsf{c}_{\tau_1} \rangle & \text{if } \rho = \tau_0 \times \tau_1, \\
\mathsf{c}_\rho x_0^{\rho_0} \cdots x_k^{\rho_k} \to \mathsf{c}_{co(\rho)} & \text{if } \rho = \rho_0 \to \cdots \to \rho_k \to co(\rho),
\end{array}
$$

with $\mathsf{c}_\iota$ and $\mathsf{c}_\epsilon$ defined to be the constants $\mathsf{c}$ and $\langle \rangle$ respectively.

2. A non-terminal $\mathsf{N}_{\pi'}^i$ for each sub-proof $\pi' \vdash B_0, \ldots, B_l$ of $\pi$ and for each $i \leq l$, with type

$$
\mathsf{N}_{\pi'}^i : \varsigma \to \tau_{B_0}^* \to \cdots \to \tau_{B_l}^* \to \tau_{B_i}
$$

and production rule(s) as given in Table 1, determined in each case by the final inference of $\pi'$.

3. A start symbol $\mathsf{S}_{\pi,i} : \tau_{A_i}$ for each $i \leq k$ with associated production rules

$$
\mathsf{S}_{\pi,i} \to \mathsf{N}_\pi^i \bot \mathsf{c}_{\tau_{A_0}^*} \cdots \mathsf{c}_{\tau_{A_k}^*}
$$

The *language of $\pi$* is the set $L(\pi) = \{(i, r^\circ) \mid i \leq k \text{ and } (\mathsf{S}_{\pi,i}, r) \in L(\mathscr{H}_\pi)\}$.

| Inference deriving $\pi$ | Corresponding production rule(s) |
|---|---|

$$\mathsf{ax}\colon \pi \vdash A, \bar{A} \qquad\qquad \mathsf{N}_\pi^i a x_0 x_1 \to \langle\rangle$$

$$\vee \frac{\pi_0 \vdash \Gamma, A, B}{\pi \vdash \Gamma, A \vee B} \qquad\qquad \mathsf{N}_\pi^i a \vec{x} z \to \begin{cases} \mathsf{N}_{\pi_0}^i a \vec{x} z z, & \text{if } i < m, \\ \langle\rangle, & \text{otherwise.} \end{cases}$$

$$\wedge \frac{\pi_0 \vdash \Gamma, A \quad \pi_1 \vdash \Delta, B}{\pi \vdash \Gamma, \Delta, A \wedge B} \qquad \mathsf{N}_\pi^i a \vec{x} \vec{y} z \to \begin{cases} \mathsf{N}_{\pi_0}^i a \vec{x} z, & \text{if } i < m, \\ \mathsf{N}_{\pi_1}^{i-m} a \vec{y} z, & m \le i < m+n, \\ \langle\rangle, & \text{otherwise.} \end{cases}$$

$$\forall_{\vec{\alpha}} \frac{\pi_0 \vdash \Gamma, A(\vec{\alpha}/\vec{v})}{\pi \vdash \Gamma, \forall \vec{v} A} \qquad \mathsf{N}_\pi^i a \vec{x} \langle z_0, \ldots, z_{p+1} \rangle \to \mathsf{N}_{\pi_0}^i ([\alpha_p \hookleftarrow z_p] \cdots [\alpha_0 \hookleftarrow z_0] a) \vec{x} z_{p+1}$$

$$\exists_{\vec{r}} \frac{\pi_0 \vdash \Gamma, A(\vec{r}/\vec{v})}{\pi \vdash \Gamma, \exists \vec{v} A} \qquad \mathsf{N}_\pi^i a \vec{x} z \to \begin{cases} \vec{r} \cdot a \star (\mathsf{N}_{\pi_0}^m a \vec{x}), & \text{if } i = m \text{ and } u(A) > 0, \\ \vec{r} \cdot a \star \langle\rangle, & \text{if } i = m \text{ and } u(A) = 0, \\ \mathsf{N}_{\pi_0}^i a \vec{x} (z(\mathsf{N}_\pi^m a \vec{x} z)) & \text{if } i \neq m, \end{cases}$$

$$\mathsf{cut} \frac{\pi_0 \vdash \Gamma, A \quad \pi_1 \vdash \Delta, \bar{A}}{\pi \vdash \Gamma, \Delta} \qquad \mathsf{N}_\pi^i a \vec{x} \vec{y} \to \begin{cases} \mathsf{N}_{\pi_0}^i a \vec{x}((\mathsf{N}_{\pi_1}^n a \vec{y}) \circ_A (\mathsf{N}_{\pi_0}^m a \vec{x})), & \text{if } i < m, \\ \mathsf{N}_{\pi_1}^{i-m} a \vec{y}((\mathsf{N}_{\pi_0}^m a \vec{x}) \circ_{\bar{A}} (\mathsf{N}_{\pi_1}^n a \vec{y})), & \text{if } m \le i. \end{cases}$$

$$\mathsf{w} \frac{\pi_0 \vdash \Gamma}{\pi \vdash \Gamma, A} \qquad\qquad \mathsf{N}_\pi^i a \vec{x} z \to \begin{cases} \mathsf{c}_{\tau_A}, & \text{if } i = m, \\ \mathsf{N}_{\pi_0}^i a \vec{x}, & \text{otherwise.} \end{cases}$$

$$\mathsf{c} \frac{\pi_0 \vdash \Gamma, A, A}{\pi \vdash \Gamma, A} \qquad\qquad \mathsf{N}_\pi^i a \vec{x} z \to \begin{cases} \mathsf{N}_{\pi_0}^i a \vec{x} z z, & \text{if } i < m, \\ \mathsf{N}_{\pi_0}^i a \vec{x} z z \mid \mathsf{N}_{\pi_0}^{i+1} a \vec{x} z z, & \text{if } i = m. \end{cases}$$

$$\mathsf{p} \frac{\pi_0 \vdash \Gamma, B, A, \Delta}{\pi \vdash \Gamma, A, B, \Delta} \qquad \mathsf{N}_\pi^i a \vec{x} z_0 z_1 \vec{y} \to \begin{cases} \mathsf{N}_{\pi_0}^{i+1} a \vec{x} z_1 z_0 \vec{y}, & \text{if } i = m, \\ \mathsf{N}_{\pi_0}^{i-1} a \vec{x} z_1 z_0 \vec{y}, & \text{if } i = m+1, \\ \mathsf{N}_{\pi_0}^i a \vec{x} z_1 z_0 \vec{y}, & \text{otherwise.} \end{cases}$$

$$\begin{aligned} \vec{\alpha} &= (\alpha_0, \ldots, \alpha_p) \\ \vec{r} &= (r_0, \ldots, r_p) \\ \vec{r} \cdot a &= (r_0 \cdot a, \ldots, r_p \cdot a) \\ (u_j)_{j \le q} \star t &= \langle u_0, \ldots, u_q, t \rangle \end{aligned} \qquad r \circ_A s = \begin{cases} r, & \text{if } e(A) > 0, \\ rs, & \text{if } u(A) > 0, \\ \langle\rangle, & \text{otherwise.} \end{cases}$$

Table 1: Production rules of $\mathscr{H}_\pi$. $\vec{x}$ and $\vec{y}$ are sequences of distinct variable symbols of length $m := |\Gamma|$ and $n := |\Delta|$ respectively.

It remains to check that the production rules of Herbrand schemes are well-typed. This task will be taken up later in Lemma 3.32. For now we take for granted the fact that Herbrand schemes are well-defined and continue with some basic properties of them (Lemmas 3.26 to 3.29) followed by the intended interpretation of the schemes as generating Herbrand disjunctions (Definition 3.30) and the observation that this interpretation coincides with the Herbrand set for quasi cut-free proofs (Lemma 3.31). We

start, however, with a brief explanation of some of the production rules from Table 1.

**Remark 3.25.** We comment on some of the rules from Table 1.

- Axiom. We are restricting axioms to quantifier-free formulæ only, which motivates the simple production rule given in the table. One may wish to permit axioms $\pi \vdash A_0, A_1$ where $A_1 = \bar{A}_0$ has arbitrary (prenex) complexity. These can be accommodated by the production rules

$$
\mathsf{N}_\pi^i a x_0 x_1 \to \begin{cases} x_{1-i}, & \text{if } u(A_i) = 0, \\ x_{1-i} x_i, & \text{if } u(A_i) > 0, \end{cases}
$$

  which the interested reader can check are well-typed. This definition mimics the behaviour of the Herbrand scheme for the natural proof of $A_0, A_1$ that uses only quantifier-free instances of axioms and alternate applications of $\exists$ and $\forall$ inferences. Our reason for favouring quantifier-free axioms is that, as a consequence, production rules never return their arguments as output, a fact that simplifies some technical aspects of the later analysis (specifically Lemma 5.9).

- $\wedge$ and $\vee$. As proofs involve prenex formulæ only, conjunctions and disjunctions are necessarily quantifier-free with associated type $\epsilon$, and therefore possess no computational content relevant to the construction of a Herbrand disjunction. When focusing on such formulæ, the production rule in each case returns the empty sequence.

- $\exists_{\vec{r}}$. The production rule in this case depends on both $i$ and the quantifier form of the active formula. Consider the instance of $\exists_{\vec{r}}$ given in Table 1. As $\pi$ is assumed regular, the active formula ($A$ in the table) is either quantifier-free or universally quantified. If $i$ marks the active formula (i.e. $i = m$) then the production rule for $\mathsf{N}_\pi^i$ directly outputs the witness terms provided by the proof and the current substitution (the sequence $(r_0 \cdot a, \dots, r_p \cdot a)$) as the first $p + 1$ components of a nested pair. The final component is either trivial (in case $A$ is quantifier-free) or, if $A$ is universally quantified, the continuation of the trace to the immediate sub-proof in the form of a function. If $i \neq m$, the production rule instead passes the above term to the corresponding argument.

- $\forall_{\vec{\alpha}}$. This is the only case that involves pattern matching in Herbrand schemes. Although it can be simulated by a recursion scheme without pattern-matching using projection functions for pair types, doing so introduces a duplication of arguments that is avoided in the chosen formulation. For instance, the production rule for $\forall_{\vec{\alpha}}$ where $\vec{\alpha}$ consists of the single eigenvariable $\alpha$ and the sequent $\Gamma$ is empty yields the production rule

$$
\mathsf{N}_\pi^0 a \langle z_0, z_1 \rangle \to \mathsf{N}_{\pi_0}^0 ([\alpha \hookleftarrow z_0] a) z_1
$$

  which may be simulated by the rule

$$
\mathsf{N}_\pi^0 a z \to \mathsf{N}_{\pi_0}^0 ([\alpha \hookleftarrow \mathsf{p}_0 z] a)(\mathsf{p}_1 z) \tag{1}
$$

where $\mathsf{p}_0$ and $\mathsf{p}_1$ are constants representing the two projection functions for pair types. If $s$ is a term such that $s \rightarrow^*_{\mathscr{H}} \langle r_0, s_0 \rangle \mid \langle r_1, s_1 \rangle$ and the four sub-terms are pairwise distinct then the reduction in (1) permits the derivation $\mathsf{N}^0_\pi \bot s \rightarrow^*$ $\mathsf{N}^0_{\pi_0}([\alpha \hookleftarrow \mathsf{p}_0 \langle r_0, s_0 \rangle] \bot)(\mathsf{p}_1 \langle r_1, s_1 \rangle)$, essentially the term $\mathsf{N}^0_{\pi_0}([\alpha \hookleftarrow r_0] \bot)s_1$, which is forbidden in the Herbrand scheme due to pattern-matching. In this sense pattern matching plays a role analogous to the rigidity conditions utilised in [24, 1, 2] for representing first-order proofs with $\Pi_1/\Pi_2$ cut complexity.

- cut. For each choice of $i$, the rule provides exactly one reduction for the non-terminal $\mathsf{N}^i_\pi$: for $i < m$ this is

$$\mathsf{N}^i_\pi a \vec{x} \vec{y} \rightarrow \begin{cases} \mathsf{N}^i_{\pi_0} a \vec{x}(\mathsf{N}^n_{\pi_1} a \vec{y}(\mathsf{N}^m_{\pi_0} a \vec{x})), & \text{if } u(A) > 0, \\ \mathsf{N}^i_{\pi_0} a \vec{x}(\mathsf{N}^n_{\pi_1} a \vec{y}), & \text{if } e(A) > 0, \\ \mathsf{N}^i_{\pi_0} a \vec{x}\langle\rangle, & \text{if } A \text{ is q.f.} \end{cases}$$

Note that, in the case $e(A) = 0$ the type $\tau^*_A$ (which marks the final argument to $\mathsf{N}^i_{\pi_0}$) is prime, and is otherwise the function type $\tau^*_{\bar{A}} \rightarrow \tau_{\bar{A}}$. Moreover, the case distinction above is independent of $i$. For instance, if $A = \forall v B$ exactly the following production rules arise from the cut.

$$\mathsf{N}^j_\pi a \vec{x} \vec{y} \rightarrow \begin{cases} \mathsf{N}^j_{\pi_0} a \vec{x}(\mathsf{N}^n_{\pi_1} a \vec{y}(\mathsf{N}^m_{\pi_0} a \vec{x})), & \text{if } j < m, \\ \mathsf{N}^{m-j}_{\pi_1} a \vec{y}(\mathsf{N}^m_{\pi_0} a \vec{x}), & \text{if } m \le j < m + n. \end{cases}$$

In the following let $\mathscr{H} = \langle \Sigma, \mathcal{N}, \mathcal{S}, \mathcal{P} \rangle$ be the Herbrand scheme for a regular proof $\pi$ with prenex $\Sigma_1$ end-sequent.

**Lemma 3.26.** *$\mathscr{H}$ is an acyclic recursion scheme. Hence, $L(\pi)$ is finite.*

*Proof.* Let $<$ be the transitive relation on non-terminals in $\mathscr{H}$ generated by the equations: $\mathsf{c}_\rho < \mathsf{c}_\sigma$ if $\rho$ is a proper sub-type of $\sigma$; $\mathsf{c}_\rho < \mathsf{N}^i_{\pi_0}$ for every $\rho$, sub-proof $\pi_0$ of $\pi$ and $i$; $\mathsf{N}^i_{\pi_0} < \mathsf{N}^j_{\pi_1}$ if either $\pi_0$ is a proper sub-proof of $\pi_1$ or $\pi_0 = \pi_1$ and $j < i$; and $\mathsf{N}^i_\pi < \mathsf{S}_{\pi,i}$ for any $i$. Clearly $<$ is acyclic and irreflexive. Moreover, for every production rule $\mathsf{F}\vec{x} \rightarrow_{\mathscr{H}} t$ and any non-terminal $\mathsf{G}$ occurring in $t$ we have $\mathsf{G} < \mathsf{F}$. $\square$

**Lemma 3.27.** *Every $\mathscr{H}$-term of simple type is a $\Sigma$-term, and every $\mathscr{H}$-term of substitution stack type has the form either $\bot$ or $[\alpha \hookleftarrow s]b$ for some $\alpha \in \Sigma$, $\Sigma$-term $s$ and $b : \varsigma$.*

*Proof.* The non-terminals of $\mathscr{H}$ all have type one of three forms: $\epsilon$, pair type, or function type with non-simple co-domain. It therefore follows that the only $\mathscr{H}$-terms of simple type are the $\Sigma$-terms. Likewise, $\bot$ and $[\alpha \hookleftarrow s]b$ are the only kind of $\mathscr{H}$-terms of type $\varsigma$. Given a substitution stack $[\alpha \hookleftarrow s]b$ however, as $\alpha \in \Sigma$ the first part of the lemma implies that $s$ is a $\Sigma$-term. $\square$

**Lemma 3.28.** *If $r : \sigma \rightarrow \tau$ is a $\mathscr{H}$-term then $\tau$ is a basic type and $\sigma$ is either basic or the type of substitution stacks. In the latter case, $r = \mathsf{N}^i_\pi$ or $\hat{\mathsf{N}}^i_\pi$ for some $\pi$ and $i$.*

*Proof.* By inspection of the types of non-terminals and terms. $\square$

**Lemma 3.29.** *Suppose $r$ is an $\mathscr{H}$-term of type $\epsilon$ containing no explicit substitutions (i.e. having no sub-term of the form $t \cdot a$). If $r \to_{\mathscr{H}}^* s$ for some basic term $s$ then $s = \langle \rangle$.*

*Proof.* By induction on the proof generating $\mathscr{H}$, on the composition of $r$ and the length of the derivation $r \to^* s$. $\qquad\square$

We now describe how Herbrand schemes can be interpreted as ascribing existential content to first-order proofs.

**Definition 3.30** (Herbrand expansion)**.** Let $\pi \vdash \Gamma$ be a proof with $\Gamma = \exists \vec{v}_0 A_0, \ldots, \exists \vec{v}_k A_k$ where $A_i$ is quantifier-free for each $i \leq k$. Let $k_i$ be the length of $\vec{v}_i$. The *Herbrand expansion of $\pi$* is the quantifier free sequent $\Gamma^\pi$ given by

$$\Gamma^\pi := \{A_i(\vec{r}_i/\vec{v}_i) \mid \vec{r}_i = (r_j)_{j<k_i} \text{ and } (i, \langle r_0, \ldots, r_{k_i-1}, \langle \rangle \rangle) \in L(\pi)\}.$$

**Lemma 3.31.** *If $\pi \vdash \Gamma$ is a quasi cut-free proof of a $\Sigma_1$ end-sequent then the Herbrand expansion of $\pi$ is a valid sequent and $\bigvee \Gamma^\pi$ is a Herbrand disjunction in the sense of Theorem 2.2.*

*Proof.* Observe that in every production rule associated to a quantifier-free cut, the term $r \circ_A s$ becomes $\langle \rangle$. Derivations in $\pi$ are therefore in 1-1 correspondence with traces following the breakdown of formulæ in the end-sequent. As a result we observe that $L(\pi)$ simply outputs all literal witnesses to the existential quantifiers in the end-sequent. $\quad\square$

The idea behind Herbrand schemes is to provide a generalisation of the above lemma to proofs containing quantified cuts. The analysis necessary for the result is carried out in Section 5. In the remainder of this section we prove the production rules of Herbrand schemes are well-typed and derive upper bounds on the size of Herbrand expansions.

**Lemma 3.32.** *The production rules of Herbrand schemes are type preserving.*

*Proof.* Fix a proof $\pi$ with prenex end-sequent $A_0, \ldots, A_m$ and $i \leq m$. We establish type-preservation of the production rules for the non-terminals $\mathsf{N}_\pi^0, \ldots, \mathsf{N}_\pi^m$ via a case distinction on the final inference rule in $\pi$.

Suppose $\pi \vdash A_0, \ldots, A_{m-1}, \exists \vec{v} A$ is obtained from proof $\pi_0$ by $\exists_{\vec{r}}$. Thus $\pi_0 \vdash \Gamma, A(\vec{r}/\vec{v})$ for some sequence $\vec{r} = (r_j)_{j \leq k}$ of simple $\Sigma$-terms of type $\iota$. By regularity, $e(A) = 0$, i.e. either $A$ is quantifier-free or $u(A) > 0$. Let $\Gamma = A_0, \ldots, A_{m-1}$ and fix a term $z : \tau_{\exists \vec{v} A}^*$ and a sequence of terms $\vec{x}$ of length $m$ such that $\mathsf{N}_\pi^i \vec{x} z$ is well-typed. By definition $\mathsf{N}_{\pi_0}^i$ has type

$$\mathsf{N}_{\pi_0}^i : \varsigma \to \tau_{A_0}^* \to \cdots \to \tau_{A_{m-1}}^* \to \tau_A^* \to \begin{cases} \tau_A, & \text{if } i = m, \\ \tau_{A_i}, & \text{otherwise.} \end{cases}$$

To check type preservation there are two cases to consider:

1. $i = m$. If $u(A) = 0$ then $A$ is quantifier-free and $\tau_A^* = \epsilon$. If $u(A) > 0$ then $\tau_A^* = \tau_A^* \to \tau_A$ by Lemma 3.22(3), so the type of $\mathsf{N}_{\pi_0}^i a \vec{x}$ is $\tau_A^*$. Since also $\tau_{\exists \vec{v} A} = \underbrace{\iota \times \cdots \times \iota}_{k} \times \tau_A^*$ by Lemma 3.22(4), we are done.

| $A$ | $\Sigma_0$ | $\Sigma_n \setminus \Pi_n$ | $\Pi_n \setminus \Sigma_n$ |
|---|---|---|---|
| $ord(\tau_A)$ | $0$ | $n \mathbin{\dot-} 2$ | $n \mathbin{\dot-} 3$ |
| $ord(\tau_A^*)$ | $0$ | $n \mathbin{\dot-} 1$ | $n \mathbin{\dot-} 2$ |

Table 2: Order of types $\tau_A$ and $\tau_A^*$.

2. $i \neq m$. In this case it is necessary to check that $\tau_{\exists \vec{v} A}^* = \tau_{\exists \vec{v} A} \to \tau_A^*$. But this follows directly from the definition and the fact that $\tau_{\forall \vec{v} \bar{A}} = \tau_{\bar{A}} = \tau_A^*$ as $e(A) = 0$.

Suppose $\pi$ is derived from $\pi_0$ via the inference $\forall_{\vec{\alpha}}$ and $A_m = \forall \vec{v} A$ with $u(A) = 0$ and $\vec{\alpha} = (\alpha_j)_{j<k}$. Let $i \leq m$ and fix terms $\vec{x}$, $\vec{z} = (z_0, \ldots, z_k)$ such that $\mathsf{N}_\pi^i a \vec{x} \langle z_0, \ldots, z_k \rangle$ is well-typed. Lemma 3.22 implies that $z_j : \iota$ for each $j < k$, and $z_k : \tau_A^*$. Thus $\mathsf{N}_{\pi_0}^i b \vec{x} z_k$ is well-typed and has type $\tau_A = \tau_{A_m}$.

Suppose $\pi$ is derived via $\mathsf{cut}$ from sub-proofs $\pi_0 \vdash \Gamma, A$ and $\pi_1 \vdash \Delta, \bar{A}$. Let $m = |\Gamma|$ and $n = |\Delta|$ and fix $\vec{x}$ and $\vec{y}$ suitably typed. Without loss of generality we may assume $i < m$, in which case we require to show $(\mathsf{N}_{\pi_1}^n a \vec{y}) \circ_A (\mathsf{N}_{\pi_0}^m a \vec{x}) : \tau_A^*$ which reduces (via Remark 3.25) to proving

$$e(A) > 0 \text{ implies } \tau_A^* = \tau_{\bar{A}}^* \to \tau_{\bar{A}},$$
$$u(A) > 0 \text{ implies } \tau_A^* = \tau_{\bar{A}} \text{ and } \tau_{\bar{A}}^* = \tau_A^* \to \tau_A,$$

both of which follow directly from Lemma 3.22.

The remaining cases are straightforward and omitted. $\qquad\square$

**Lemma 3.33.** *Fix a prenex formula $A$. The order of $\tau_A$, $\tau_A^*$ are as presented in Table 2 where $\dot-$ denotes subtraction truncated at 0, i.e. $n \mathbin{\dot-} m = \max\{n - m, 0\}$.*

*Proof.* By induction on complexity of $A$. If $A$ is quantifier free then $\tau_A = \epsilon = \tau_A^*$ so $ord(\tau_A) = ord(\tau_A^*) = 0$. Moreover, by Example 3.23, the lemma holds for $A \in (\Sigma_1 \setminus \Pi_1) \cup (\Pi_1 \setminus \Sigma_1)$. Suppose $n > 1$. For $A = \exists \vec{v} B$ where $B \in \Pi_{n-1} \setminus \Sigma_{n-1}$,

$$\begin{aligned}
ord(\tau_A) &= ord(\tau_{\bar{B}}^*) & \text{(Lemma 3.22(4))} \\
&= n - 2 & \text{(induction hypothesis)} \\
ord(\tau_A^*) &= \max\{ord(\tau_A) + 1, ord(\tau_{\bar{B}})\} & \text{(definition)} \\
&= n - 1 & \text{(induction hypothesis)}
\end{aligned}$$

For $A = \forall \vec{v} B$ where $B \in \Sigma_{n-1} \setminus \Pi_{n-1}$,

$$\begin{aligned}
ord(\tau_A) &= ord(\tau_B) & \text{(definition)} \\
&= n \mathbin{\dot-} 3 & \text{(induction hypothesis)} \\
ord(\tau_A^*) &= ord(\tau_{\bar{A}}) & \text{(definition)} \\
&= n - 2 & \square
\end{aligned}$$

**Corollary 3.34.** *For a proof $\pi \vdash A_0, \ldots, A_k$ and $i < k$, the order of the non-terminal $\mathsf{N}_\pi^i$ is equal to the smallest $n$ such that $\{A_j : j \leq k\} \subset \Pi_{n+1}$, unless $A_i$ is $\Pi_1$, in which case the order of $\mathsf{N}_\pi^i$ is zero.*

*Proof.* If $A_i$ is $\Pi_1$ then $\tau_{A_i} = \epsilon$ and the order of $\mathsf{N}_\pi^i$ is 0 by definition. Otherwise, the order of $\mathsf{N}_\pi^i$ is one greater than the maximum among the orders of $\tau_{A_j}^*$ for $j \leq k$. $\qquad\square$

It is now possible to strengthen Lemma 3.26 to a concrete bound on the number of terms derivable from a Herbrand scheme. The idea is to eliminate occurrences of pattern-matching in a Herbrand scheme $\mathscr{H}$ in a way that does not decrease the length of derivations so that Theorem 3.17 and Corollary 3.18 can be applied.

**Theorem 3.35.** *If $\pi \vdash \Gamma$ is a proof of a single prenex $\Sigma_1$ formula in which all cut formulæ are contained in $\Pi_n \cup \Sigma_n$ then the size of the Herbrand expansion $\Gamma^\pi$ is no greater than $2_{n+2}^{4|\pi|^3}$ where $|\pi|$ is the number of inference rules in $\pi$.*

*Proof.* The case $n = 0$ is covered by Lemma 3.31 so suppose $n > 0$. Let $\mathscr{H}$ be the Herbrand scheme of $\pi$. Since the cut rank of $\pi$ is bounded by $n$, Corollary 3.34 implies that the order of $\mathscr{H}$ is no greater than $n$. To obtain the desired bounds we apply Theorem 3.17. However, this requires first eliminating the explicit substitutions introduced by the $\forall$ inferences. Let $\mathscr{H}'$ denote the higher order recursion scheme with non-terminals of basic type obtained from $\mathscr{H}$ by removing all substitutions terms and types from non-terminals and production rules. In particular, the productions originating from $\forall_{\vec{\alpha}}$ and $\exists_{\vec{r}}$ inferences are replaced by following in $\mathscr{H}'$:

$$\forall_{\vec{\alpha}}: \ \mathsf{N}_{\pi'}^i \vec{x} \langle z_0, \ldots, z_{p+1} \rangle \to \mathsf{N}_{\pi_0'}^i \vec{x} z_p$$

$$\exists_{\vec{r}}: \qquad\qquad \mathsf{N}_\pi^i \vec{x} z \to \begin{cases} \mathsf{N}_{\pi_0}^i \vec{x}(z(\mathsf{N}_\pi^m \vec{x})), & i \neq m, \\ \langle \mathsf{c}, \ldots, \mathsf{c}, \mathsf{N}_{\pi_0}^m \vec{x} \rangle, & i = m \text{ and } u(A) > 0, \\ \langle \mathsf{c}, \ldots, \mathsf{c}, \langle \rangle \rangle, & i = m \text{ and } u(A) = 0. \end{cases}$$

The second part of Lemma 3.27 implies that derivations in $\mathscr{H}'$ from the start symbol are in 1-1 correspondence with derivations in $\mathscr{H}$. Repeating the argument of Corollary 3.18, the size of $L(\pi)$ is therefore bounded by $2^K$ where $K$ is the length of the longest derivation in $\mathscr{H}'$ from the single start symbol. The order of $\mathscr{H}'$ is no greater than $n$, the number of non-terminals is bounded by $|\pi|^2$, and for each production rule $\mathsf{F}\vec{x} \to t$ in $\mathscr{H}'$, $|t|_\Sigma < 3 \times |\pi|$ where $\Sigma$ is the ranked alphabet of function symbols and constants occurring in $\pi$. Theorem 3.17 then implies $K \leq 2_{n+1}^{4|\pi|^3}$. $\qquad\square$

# 4 Example: a Herbrand Disjunction for the Pigeonhole Principle

We consider a formal proof of the pigeonhole principle for two boxes via the infinite pigeonhole principle. The question of the computational content of this proof is attributed to G. Stolzenberg in [12]. A variety of analytic methods have since been applied to this proof [23, 9, 43, 7, 1] and its generalisations [40, 38]. The version we present here is a formal proof with a single $\Pi_3$ cut based on the proof with two $\Pi_2$ cuts given in [1, 43].

Let $f: \mathbb{N} \to \{0, 1\}$ be a total Boolean function, let $I_i$ (for $i = 0, 1$) express that there are infinitely many $m \in \mathbb{N}$ for which $f(m) = i$ and $T$ express that there exists $m < n$ such that $f(m) = f(n)$. A consequence of the law of excluded middle is $\exists w I_w$. Moreover,

$I_i$ implies $T$ for each $i \in \{0,1\}$: assuming $I_i$ there exists $m \geq 0$ and $n \geq m+1$ for which $f(m) = f(n) = i$. Combining these observations we conclude $T$.

The following formalises the above argument into a proof with a single $\Pi_3$ cut. The formal language, $\Sigma$, comprises two unary function symbols $\mathsf{f}$, $\mathsf{s}$, one binary function symbol $\mathsf{m}$, a constant symbol $0$ and a binary relation $\leq$. We make the following definitions and abbreviations:

- $T = \exists u \exists v (u < v \wedge \mathsf{f}u = \mathsf{f}v)$,

- $I = \exists w I_w$ where $I_r = \forall u \exists v (u \leq v \wedge \mathsf{f}v = r)$,

- $\Gamma = \{\forall u \forall v (u \leq \mathsf{m}uv \wedge v \leq \mathsf{m}uv), \forall u (\mathsf{f}u = 0 \vee \mathsf{f}u = \mathsf{s}0)\}$,

- $\Delta = \{\forall u \forall v \forall w (u = v \wedge w = v \rightarrow u = w), \forall u \forall v (\mathsf{s}u \leq v \rightarrow u < v)\}$,

- $I_r^s$ and $I_r^{s,t}$ denote, respectively, $\exists v (s \leq v \wedge \mathsf{f}v = r)$ and $(s \leq t \wedge \mathsf{f}t = r)$,

- $T_{s,t}$ denotes $(s < t \wedge \mathsf{f}s = \mathsf{f}t)$.

The intended interpretation of the symbols is: $\mathsf{f}$ represents the (arbitrary) function $f$, $\mathsf{s}$ the successor function on $\mathbb{N}$, $\leq$ the standard ordering and $\mathsf{m}$ the binary max function.



Figure 4: Proof $\pi_\infty$ of pigeonhole principle.

A formal proof of the pigeonhole principle (namely $\Gamma, \Delta \vdash T$) is given in Figure 4 which we name $\pi_\infty$. The proof is displayed in two-sided sequent calculus as this simplifies the presentation and following discussion. The intended interpretation of the two-sided sequent $A_1, \ldots, A_k \vdash B_1, \ldots, B_l$ is the sequent $\bar{A}_1, \ldots, \bar{A}_k, B_1, \ldots, B_l$. For brevity, only eigenvariables and witnesses of the quantifiers and instances of the existential formula $T$ are displayed in $\pi_\infty$. The proof fully fleshed out uses about 50 application of the axioms and rules of the calculus but the only cut in $\pi_\infty$ is the one displayed in the figure. Two normal forms of the proof of size $\sim 200$ have been computed in a case study [43] from which one can read off the Herbrand sets for the formula $T$ (also for formulæ in $\Gamma \cup \Delta$ but these are less interesting). Up to interpretation of the logical symbols by their intended semantics, the two Herbrand sets combined provide

the witnesses $\{\langle 0,1\rangle, \langle 1,2\rangle, \langle 2,3\rangle, \langle 0,2\rangle, \langle 1,3\rangle\}$ to the existential quantifiers in $T$.[1] The Herbrand scheme $\mathscr{H}_{\pi_\infty}$ associated to the proof $\pi_\infty$ computes the same Herbrand set, a fact we demonstrate in the following.

**Types and terms**  The Herbrand scheme for $\pi_\infty$ comprises a non-terminal for each sub-proof of $\pi_\infty$ and each formula in the end-sequent of that sub-proof. Recall, for each sub-proof $p : \Pi \vdash \Lambda$ of $\pi_\infty$ and each $i < |\Pi| + |\Lambda|$ there is a non-terminal $\mathsf{N}_p^i$ in $\mathscr{H}_{\pi_\infty}$ representing the existential content of the $i$-th formula in the sequent at position $p$. In the following, in place of $\mathsf{N}_p^i$ we will write $\mathsf{N}_p^A$ where $A$ is the $i$-th formula in the sequent assuming this is unique. In case $A$ occurs more than once in the sequent $\Pi \vdash \Lambda$ (such as at positions $b$ and 2) the non-terminal $\mathsf{N}_p^A$ refers to the first occurrence of $A$ and we use the notation $\mathsf{N}_p^{A^+}$ for the second occurrence. Concerning the type of $\mathsf{N}_p^A$, we recall the types $\tau_F$ and $\tau_F^*$ for each formula $F$ in $\pi_\infty$. Let $\hat{\iota} = \iota \times (\iota^1 \to \epsilon)$ and $\hat{\epsilon} = \iota^2 \to \epsilon$.

- For $F \in \Gamma \cup \Delta$ we have $\tau_{\bar{F}} \in \{\iota^1, \iota^2, \iota^3\}$, and $\tau_{\bar{F}}^* = \tau_{\bar{F}} \to \epsilon$.

- $T$: $\tau_T = \iota^2$, $\tau_T^* = \hat{\epsilon}$.

- $I_r^s$: $\tau_{I_r^s} = \iota^1 = \tau_{\bar{I}_r^s}^*$, $\tau_{I_r^s}^* = \iota^1 \to \epsilon$ and $\tau_{\bar{I}_r^s} = \epsilon$.

- $I_r$: $\tau_{I_r} = \iota^1$, $\tau_{\bar{I}_r} = \tau_{I_r}^* = \hat{\iota}$ and $\tau_{\bar{I}_r}^* = \hat{\iota} \to \iota^1$.

- $I$: $\tau_I = \iota \times (\hat{\iota} \to \iota^1) = \tau_{\bar{I}}^*$, $\tau_{\bar{I}} = \hat{\iota}$ and $\tau_I^* = \tau_I \to \tau_{\bar{I}}$.

- The remaining formulæ that occur in $\pi_\infty$ are quantifier-free and are assigned type $\epsilon$ in all cases.

According to the definition, the type of $\mathsf{N}_{\pi_\infty}^T$ is $\varsigma \to \tau_{\bar{F}}^* \to \tau_{\bar{G}}^* \to \tau_{\bar{C}}^* \to \tau_{\bar{D}}^* \to \tau_T^* \to \tau_T$ where $\varsigma$ is the type of substitution stacks, $F$ and $G$ are the two formulæ in $\Gamma$ and $C$ and $D$ are the formulæ in $\Delta$. As the formulæ in $\Gamma \cup \Delta$ are $\Sigma_1$, their input type carries no computational content (cf. Lemma 5.17), and we can ignore these formulæ and identify the type above with $\varsigma \to \tau_T^* \to \tau_T$, and the term $\mathsf{N}_{\pi_\infty}^T a \mathsf{c}_{\tau_{\bar{F}}^*} \mathsf{c}_{\tau_{\bar{G}}^*} \mathsf{c}_{\tau_{\bar{C}}^*} \mathsf{c}_{\tau_{\bar{D}}^*}$ with $\mathsf{N}_{\pi_\infty}^T a$. Likewise, the type of $\mathsf{N}_c^{I_1}$ is assumed to be $\varsigma \to \tau_I^* \to \tau_{I_1}^* \to \tau_{I_1}$ and the type of $\mathsf{N}_2^{\bar{I}_\gamma^+}$ is $\varsigma \to (\hat{\iota} \to \iota^1) \to (\hat{\iota} \to \iota^1) \to \hat{\epsilon} \to \hat{\iota}$.

Other abbreviations and simplifications we utilise are:

- $\langle r \rangle$ for either the sequence $\langle r, \langle\rangle\rangle$ or $\langle r, \mathsf{c}_{\iota^1 \to \epsilon}\rangle$, depending on type, and $\langle r, s \rangle$ as a term of type $\iota^2$ represents $\langle r, s, \langle\rangle\rangle$.

- $\hat{0} = \mathsf{m}00$, $1 = \mathsf{s}\hat{0}$, $\hat{1} = \mathsf{m}01$, $2 = \mathsf{s}(\mathsf{m}10)$ and $\hat{2} = \mathsf{m}02$.

- For each non-terminal $\mathsf{N}_p^A$ where $A$ is the $i$-th formula at position $p$, an additional non-terminal $\hat{\mathsf{N}}_p^A$ with the same arity as $\mathsf{N}_p^A$ and associated production rule

$$\hat{\mathsf{N}}_p^A a x_0 \cdots x_k \to \mathsf{N}_p^A a x_0 \cdots x_{i-1} x_k x_i \cdots x_{k-1}$$

is included in the Herbrand scheme $\mathscr{H}_{\pi_\infty}$. These non-terminals ease the computation in derivation steps involving permutation.

---

[1] In [43] $\pi_\infty$ is formalised as a proof with two $\Pi_2$-cuts but as far as computing Herbrand sets, the two proofs are essentially identical.

- The Herbrand scheme also includes explicit non-terminals for non-determinism at each type, which are represented via set notation: for terms $s_0, \ldots, s_k : \rho$ of the same type, the set $S = \{s_i \mid i \leq k\}$ is a term of type $\rho$ with reduction $S \to s_i$ for each $i \leq k$.

- An equivalence relation $\asymp$ on terms of identical type defined as inducing the same language within all contexts. Formally, we set $r \asymp s$ iff $r \prec s \prec r$ where $r \prec s$ holds just if $r, s : \rho$ and for every $\mathscr{H}_{\pi_\infty} \cup \{x^\rho\}$-term $t$ of basic type (where $x$ is a fresh symbol of type $\rho$), whenever $t(r/x) \to^* u$ for a $\Sigma$-term $u$, then $t(s/x) \to^* v$ for some $\Sigma$-term $v$ such that $u^\circ = v^\circ$.

For instance, if $r \to s$ via an application of a deterministic production rule then $r \asymp s$, and if $S = S'$ are two representations of the same set of terms then $S \asymp S'$. In general, $r(S/x) \not\asymp \{r(s/x) \mid s \in S\}$ as shown by considering $r = \mathsf{F}x$ with reduction $\mathsf{F}x \to^* \mathsf{m}xx$.[2] However, suppose $r = \mathsf{F}t_1 \cdots t_k x$, $S$ is a set of pairs, $\mathsf{F}$ is deterministic and $\mathsf{F}t_1 \cdots t_k \langle x, x' \rangle \to t$. Then $r(S/x) \asymp \{r(s/x) \mid s \in S\} \asymp \{t((u,v)/(x,x')) \mid \langle u, v \rangle \in S\}$.

Finally, we remark that, generalising Lemma 3.29, for every type $\rho$ with co-domain $\epsilon$ and every term $r : \rho$, we have $r \asymp \mathsf{c}_\rho$.

**Language of $\pi_\infty$**   We now compute the language of $\mathscr{H}_{\pi_\infty}$ focussing on the formula $T$, i.e. set of terms (after evaluation) derivable from the term $\mathsf{N}_{\pi_\infty}^T \bot \mathsf{c}_{\hat{\epsilon}}$. The first, and only, production rule applicable to this term is given by the cut rule at the root of the proof:

$$\mathsf{N}_{\pi_\infty}^T \bot \mathsf{c}_{\hat{\epsilon}} \to \mathsf{N}_0^T \bot (\mathsf{N}_a^I \bot (\hat{\mathsf{N}}_0^{\bar{I}} \bot \mathsf{c}_{\hat{\epsilon}})) \mathsf{c}_{\hat{\epsilon}}. \tag{2}$$

Analysing derivations directly from this term is complicated. As the right sub-proof at $0$ culminates in a $\forall_\gamma$ inference, the external non-terminal $\mathsf{N}_0^{\bar{I}}$ cannot be reduced until its second argument (the term $\mathsf{N}_a^I \bot (\hat{\mathsf{N}}_0^{\bar{I}} \bot \mathsf{c}_{\hat{\epsilon}})$) is reduced to an explicit pair. But the inference at $a$ in the left sub-proof is a contraction, so this immediately introduces non-determinism and duplication of arguments. After resolving the non-determinism and reducing the two continuations of $\mathsf{N}_a^I$ to pairs (say in terms of $\mathsf{N}_d^{I_0}/\mathsf{N}_d^{I_1}$), the external non-terminal can be reduced. The argument $\hat{\mathsf{N}}_0^{\bar{I}} \bot \mathsf{c}_{\hat{\epsilon}}$ comes into play at this point: the productions for $\mathsf{N}_b^{I^+}$ and $\mathsf{N}_c^I$ increase the nesting of non-terminals which must also be evaluated as pairs in order to proceed beyond $\mathsf{N}_d^{I_0}/\mathsf{N}_d^{I_1}$.

In the following, we compute the language via a top-down approach, analysing derivations starting from relatively simple terms, and building these together to compute the language of more complex interactions between non-terminals. We begin with the most simple derivations available. Recall that $\tau_{I_0}^* = \tau_{I_1}^* = \iota^1$. Concerning non-terminals from the left sub-proof we have the following derivation starting from $\mathsf{N}_d^{I_0}/\mathsf{N}_d^{I_1}$.

$$\begin{aligned}
\mathsf{N}_d^{I_i} a \langle r \rangle \langle s \rangle &\to \mathsf{N}_e^{2+i} ([\hat{\alpha} \hookleftarrow s]a) \langle r \rangle \langle \rangle \\
&\to \mathsf{N}_f^{2+i} ([\alpha \hookleftarrow r][\hat{\alpha} \hookleftarrow s]a) \langle \rangle \langle \rangle \\
&\to^* \langle \mathsf{m}\alpha\hat{\alpha} \cdot [\alpha \hookleftarrow r][\hat{\alpha} \hookleftarrow s]a \rangle
\end{aligned}$$

---

[2] Formally, we require an $\mathscr{H}_{\pi_\infty}$ analogue of $\mathsf{F}$ but this is not difficult to find.

(Note, as $|\Gamma| = 2$ the $(2+i)$-th formula at positions $e$ and $f$ is the ancestor of $I_i$ from $d$.)
If $r$ happens to be such that $r \cdot [\hat{\alpha} \leftarrow s]a \asymp r \cdot a$, then since the derivation above follows deterministic reductions only, we deduce

$$\mathsf{N}_d^{I_i} a \langle r \rangle \langle s \rangle \asymp \langle \mathsf{m} r s \cdot a \rangle. \tag{3}$$

Examining the non-terminals from lower in the left sub-proof affords us

$$\mathsf{N}_c^I ars \to \langle 0, \hat{\mathsf{N}}_d^{I_0} as \rangle \qquad\qquad \mathsf{N}_c^{I_1} ars \to^* \mathsf{N}_d^{I_1} a(r \langle 0, \hat{\mathsf{N}}_d^{I_0} as \rangle)s$$

$$\mathsf{N}_b^{I^+} ars \to \langle 1, \mathsf{N}_c^{I_1} ar \rangle \qquad\qquad \mathsf{N}_b^I ars \to \mathsf{N}_c^I ar(s(\mathsf{N}_b^{I^+} ars))$$

$$\to^* \langle 0, \hat{\mathsf{N}}_d^{I_0} a(s \langle 1, \mathsf{N}_c^{I_1} ar \rangle) \rangle$$

The derivation from $\mathsf{N}_c^{I_1} ars$ can be continued provided that the two arguments of $\mathsf{N}_d^{I_1}$, namely $r \langle 0, \hat{\mathsf{N}}_d^{I_0} as \rangle$ and $s$, are reducible to pairs. Thus if $r \langle 0, \hat{\mathsf{N}}_d^{I_0} a \langle s \rangle \rangle \to^* \langle r_0, r_0' \rangle$ and $r_0 \cdot [\hat{\alpha} \leftarrow s]a \asymp r_0 \cdot a$ then

$$\mathsf{N}_c^{I_1} ar \langle s \rangle \to^* \mathsf{N}_d^{I_1} a \langle r_0 \rangle \langle s \rangle$$

$$\asymp \langle \mathsf{m} r_0 s \cdot a \rangle.$$

Because the reductions governing $\mathsf{N}_c^{I_1}$ are all deterministic, when phrased in terms of equivalences, this becomes

$$\left. \begin{array}{l} r \langle 0, \hat{\mathsf{N}}_d^{I_0} a \langle s \rangle \rangle \asymp \{ \langle r_i \rangle \mid i \le k \} \\ r_i \cdot [\hat{\alpha} \leftarrow s]a \asymp r_i \cdot a \text{ each } i \le k \end{array} \right\} \text{ implies } \mathsf{N}_c^{I_1} ar \langle s \rangle \asymp \{ \langle \mathsf{m} r_i s \cdot a \rangle \mid i \le k \} \tag{4}$$

Property (4) will be useful later.

Returning briefly to the derivations from non-terminals $\mathsf{N}_b^I$ and $\mathsf{N}_b^{I^+}$ started earlier, each of these derivations is also deterministic, so therefore

$$\mathsf{N}_a^I ar \asymp \{ \mathsf{N}_b^I arr, \mathsf{N}_b^{I^+} arr \}$$

$$\asymp \{ \langle 0, \hat{\mathsf{N}}_d^{I_0} a(r \langle 1, \mathsf{N}_c^{I_1} ar \rangle) \rangle, \langle 1, \mathsf{N}_c^{I_1} ar \rangle \} \tag{5}$$

which provides the first step in the continuation of the derivation from $\mathsf{N}_{\pi_\infty}^T$. Before extending (2) however we consider some simple derivations arising from the right sub-proof. On this side, the alternation of universal and existential inference rules means that few non-terminals can be adequately analysed in isolation as we did above. Most straightforward are non-terminals $\mathsf{N}_4^{\bar{I}\gamma}$ and $\mathsf{N}_2^{\bar{I}\gamma}$, for which we have

$$\mathsf{N}_4^{\bar{I}\gamma} arst \asymp \langle \mathsf{s}\beta \cdot a \rangle \qquad \mathsf{N}_2^{\bar{I}\gamma} arst \asymp \langle 0 \cdot a \rangle \asymp \langle 0 \rangle$$

This gives rise to, for example,

$$\mathsf{N}_3^{\bar{I}\gamma} a \langle r_0 \rangle st \asymp \mathsf{N}_4^{\bar{I}\gamma} ([\beta \leftarrow r_0]a) \langle \rangle st \asymp \langle \mathsf{s} r_0 \cdot a \rangle$$

and hence if $r : (\iota \times (\iota^1 \to \epsilon)) \to \iota^1$ is a term such that $r \langle 0 \rangle \asymp \{ \langle r_i \rangle \mid i \le k \}$ then also

$$\mathsf{N}_2^{\bar{I}\gamma^+} arst \asymp \mathsf{N}_3^{\bar{I}\gamma} a(r(\mathsf{N}_2^{\bar{I}\gamma} arst))st$$

$$\asymp \mathsf{N}_3^{\bar{I}\gamma} a \{ \langle r_i \rangle \mid i \le k \} st$$

$$\asymp \{ \langle \mathsf{s} r_i \cdot a \rangle \mid i \le k \}$$

The equivalences for $\mathsf{N}_2^{\bar{I}\gamma}$ and $\mathsf{N}_2^{\bar{I}\gamma^+}$ combine to yield, given the same $r$,

$$\mathsf{N}_1^{\bar{I}\gamma}art \asymp \{\mathsf{N}_2^{\bar{I}\gamma}arrt, \mathsf{N}_2^{\bar{I}\gamma^+}arrt\} \asymp \{\langle 0 \rangle\} \cup \{\langle \mathsf{s}r_i \cdot a \rangle \mid i \leq k\}. \tag{6}$$

In particular, choosing $r = \hat{\mathsf{N}}_d^{I_0} \perp \langle s \rangle$, this implies

$$\mathsf{N}_1^{\bar{I}\gamma}a(\hat{\mathsf{N}}_d^{I_0} \perp \langle s \rangle)t \asymp \{\langle 0 \rangle, \langle \mathsf{s}(\mathsf{m}0s) \cdot a \rangle\} \tag{7}$$

which will be needed later. In addition to (7), it is necessary to analyse the complex term $\mathsf{N}_1^{\bar{I}\gamma}a(\mathsf{N}_c^{I_1} \perp (\hat{\mathsf{N}}_0^{\bar{I}} \perp \mathsf{c}_{\hat{\epsilon}}))\mathsf{c}_{\hat{\epsilon}}$. However, here we can use (6) again. If $\delta : \iota$ is a fresh symbol then, applying (7) and (4) (using $r = \mathsf{N}_0^{\bar{I}} \perp \mathsf{c}_{\hat{\epsilon}}$), we get

$$\mathsf{N}_0^{\bar{I}} \perp \langle 0, \hat{\mathsf{N}}_d^{I_0} \perp \langle \delta \rangle \rangle \mathsf{c}_{\hat{\epsilon}} \asymp \mathsf{N}_1^{\bar{I}\gamma}([\gamma \leftarrow 0] \perp)(\hat{\mathsf{N}}_d^{I_0} \perp \langle \delta \rangle)\mathsf{c}_{\hat{\epsilon}}$$
$$\asymp \{\langle 0 \rangle, \langle \mathsf{s}(\mathsf{m}0\delta) \rangle\} \tag{8}$$
$$\mathsf{N}_c^{I_1} \perp (\hat{\mathsf{N}}_0^{\bar{I}} \perp \mathsf{c}_{\hat{\epsilon}})\langle \delta \rangle \asymp \{\langle \mathsf{m}0\delta \rangle, \langle \mathsf{m}(\mathsf{s}(\mathsf{m}0\delta))\delta \rangle\} \tag{9}$$

whence (6) implies

$$\mathsf{N}_1^{\bar{I}\gamma}a(\mathsf{N}_c^{I_1} \perp (\hat{\mathsf{N}}_0^{\bar{I}} \perp \mathsf{c}_{\hat{\epsilon}}))\mathsf{c}_{\hat{\epsilon}} \asymp \{\langle 0 \rangle, \langle 1 \rangle, \langle 2 \rangle\}. \tag{10}$$

We have still not examined derivations starting from the non-terminals $\mathsf{N}_0^T$, $\mathsf{N}_1^T$, and $\mathsf{N}_i^T$ for $i \geq 2$, which will arise in the computation of $L(\pi_\infty)$. The first three non-terminals behave according to

$$\mathsf{N}_0^T a\langle r, s \rangle t \asymp \mathsf{N}_1^T([\gamma \leftarrow r]a)st$$
$$\asymp \mathsf{N}_2^T([\gamma \leftarrow r]a)sst$$
$$\asymp \mathsf{N}_3^T([\gamma \leftarrow r]a)(s\langle 0 \rangle)st$$

The remaining behave similarly to the $\mathsf{N}_i^A$ non-terminals analysed earlier, except that it is $\mathsf{N}_6^T$ that provides the only 'outputs' in the derivation. In particular,

$$\mathsf{N}_6^T ars \asymp \langle \beta, \hat{\beta} \rangle \cdot a \qquad \mathsf{N}_5^T ar\langle s_0 \rangle t \asymp \langle \beta, s_0 \rangle \cdot a$$
$$\mathsf{N}_3^T a\langle r_0 \rangle st \asymp \mathsf{N}_5^T([\beta \leftarrow r_0]a)\langle\rangle(s\langle \mathsf{s}r_0 \cdot a \rangle)t$$

Let $\delta : \iota$ be a fresh symbol. Combining the two sets of equations above, if $s : \hat{\iota} \to \iota^1$ is such that $s\langle \delta \rangle \asymp \{\langle s_i \rangle \mid i \leq k\}$ and $s_i$ contains neither $\beta$ or $\gamma$ for each $i$, it follows that

$$\mathsf{N}_0^T \perp \langle r, s \rangle t \asymp \{\langle s_i \cdot [\delta \leftarrow 0] \perp, s_j \cdot [\delta \leftarrow \mathsf{s}s_i][\delta \leftarrow 0] \perp \rangle \mid i, j \leq k\}. \tag{11}$$

We can now proceed with calculating the language of $\mathsf{N}_{\pi_\infty}^T \perp \mathsf{c}_{\hat{\epsilon}}$. Let $w = \hat{\mathsf{N}}_0^{\bar{I}} \perp \mathsf{c}_{\hat{\epsilon}}$. Following on from (2) and (5) we have

$$\mathsf{N}_{\pi_\infty}^T \perp \mathsf{c}_{\hat{\epsilon}} \asymp \mathsf{N}_0^T \perp (\mathsf{N}_a^I \perp w)\mathsf{c}_{\hat{\epsilon}}$$
$$\asymp \left\{ \mathsf{N}_0^T \perp \langle 0, \hat{\mathsf{N}}_d^{I_0} \perp (w\langle 1, \mathsf{N}_c^{I_1} \perp w \rangle) \rangle \mathsf{c}_{\hat{\epsilon}}, \mathsf{N}_0^T \perp \langle 1, \mathsf{N}_c^{I_1} \perp w \rangle \mathsf{c}_{\hat{\epsilon}} \right\} \tag{12}$$

Thus, we need only compute

$$\mathsf{N}_d^{I_0}\bot\langle\delta\rangle(w\langle 1, \mathsf{N}_c^{I_1}\bot w\rangle) \qquad\qquad \mathsf{N}_c^{I_1}\bot(\hat{\mathsf{N}}_0^{\bar{I}}\bot\mathsf{c}_{\hat{\epsilon}})\langle\delta\rangle$$

and apply (11) (assuming that the terms obtained will be free of $\beta$ and $\gamma$). The latter was already established in (9):

$$\mathsf{N}_c^{I_1}\bot(\hat{\mathsf{N}}_0^{\bar{I}}\bot\mathsf{c}_{\hat{\epsilon}})\langle\delta\rangle \asymp \{\langle\mathsf{m}0\delta\rangle, \langle\mathsf{m}(\mathsf{s}(\mathsf{m}0\delta))\delta\rangle\}$$

For the former, we have $w\langle 1, \mathsf{N}_c^{I_1}\bot w\rangle \asymp \mathsf{N}_1^{\bar{I}\gamma}([\gamma\leftarrow 1]\bot)(\mathsf{N}_c^{I_1}\bot w)\mathsf{c}_{\hat{\epsilon}}$, whence (10) implies

$$\begin{aligned}
\mathsf{N}_d^{I_0}\bot\langle\delta\rangle(w\langle 1, \mathsf{N}_c^{I_1}\bot w\rangle) &\asymp \mathsf{N}_d^{I_0}\bot\langle\delta\rangle(\mathsf{N}_1^{\bar{I}\gamma}([\gamma\leftarrow 1]\bot)(\mathsf{N}_c^{I_1}\bot w)\mathsf{c}_{\hat{\epsilon}})\\
&\asymp \left\{\mathsf{N}_d^{I_0}\bot\langle\delta\rangle\langle s\rangle \mid s\in\{0,1,2\}\right\}\\
&\asymp \{\langle\mathsf{m}\delta 0\rangle, \langle\mathsf{m}\delta 1\rangle, \langle\mathsf{m}\delta 2\rangle\}
\end{aligned}$$

Hence, by (11) and (12), we deduce

$$\begin{aligned}
\mathsf{N}_{\pi_\infty}^T\bot\mathsf{c}_{\hat{\epsilon}} \asymp &\left\{\langle r, s^{[\delta\leftarrow\mathsf{s}r]}\rangle \mid r\in\{\hat{0}, \hat{1}, \hat{2}\}, s\in\{\mathsf{m}\delta 0, \mathsf{m}\delta 1, \mathsf{m}\delta 2\}\right\}\\
&\cup \left\{\langle r, s^{[\delta\leftarrow\mathsf{s}r]}\rangle \mid r\in\{\hat{0}, \mathsf{m}10\}, s\in\{\mathsf{m}0\delta, \mathsf{m}(\mathsf{s}(\mathsf{m}0\delta))\delta\}\right\}
\end{aligned}$$

Under the standard interpretation of the symbols 0, $\mathsf{s}$ and $\mathsf{m}$ (as zero, successor and binary 'max') $L(\pi_\infty)$ ascribes to $T$ the set

$$\{(0,1), (0,2), (1,2), (2,3), (1,3)\}.$$

# 5 Language Preservation

Recall the relation $\pi \rightsquigarrow \pi'$ which expresses that $\pi'$ is obtained from $\pi$ by the application of a reduction rule in Figures 2 and 3 to a sub-proof of $\pi$. In the present section we determine in which cases $\rightsquigarrow$ supports: (i) *language inclusion*: $\pi \rightsquigarrow \pi'$ implies $L(\pi') \subseteq L(\pi)$; and (ii) *language equality*: $\pi \rightsquigarrow \pi'$ implies $L(\pi') = L(\pi)$. Establishing language inclusion for cut reduction steps will suffice to derive the main theorem; language equality allows a more fine-grained study of the Herbrand content of proofs as if $\pi_0$ and $\pi_1$ are two proofs that can be connected by a sequence of forward and backward language preserving reductions then $L(\pi_0) = L(\pi_1)$.

The first task is to define a single unified Herbrand scheme in which we can reason about derivations concerning arbitrary regular proofs. At the same time, we expand the grammar by certain non-terminals that will facilitate the analysis. These additional non-terminals were introduced informally in the example of the previous section.

**Definition 5.1** (Universal Herbrand Scheme)**.** Let $\Sigma$ be the signature of first-order logic. We let $\mathscr{H}$ denote the infinite recursion scheme comprising:

1. a non-deterministic non-terminal $\mathsf{D}_\rho : \rho \to \rho \to \rho$ for each basic type $\rho$ with production rules $\mathsf{D}_\rho rs \to r$ and $\mathsf{D}_\rho rs \to s$,

2. all non-terminals $\mathsf{N}_\pi^i$, $\mathsf{S}_{\pi,i}$ and $\mathsf{c}_\rho$ from Definition 3.24 with their associated production rules formulated deterministically in terms of the $\mathsf{D}_\rho$ non-terminals above,

3. for each non-terminal $\mathsf{N}_\pi^i : \varsigma \to \tau_0 \to \cdots \tau_m \to \tau$ from the above with $\tau$ prime, a non-terminal $\hat{\mathsf{N}}_\pi^i$ with type and associated production rule

$$\hat{\mathsf{N}}_\pi^i : \varsigma \to \tau_0 \to \cdots \to \tau_{i-1} \to \tau_{i+1} \to \cdots \tau_m \to \tau_i \to \tau$$

$$\hat{\mathsf{N}}_\pi^i a x_0 \cdots x_m \to \mathsf{N}_\pi^i a x_0 \cdots x_{i-1} x_m x_i \cdots x_{m-1}$$

We refer to $\mathscr{H}$ as the *universal Herbrand scheme*.

Henceforth, a *term* is an $\mathscr{H}$-term and we write $\to$ in place of $\to_{\mathscr{H}}$. Finite sets of $\mathscr{H}$-terms will represent applications of the non-deterministic non-terminals $\mathsf{D}_\rho$. Specifically, the set $\{s_0^\rho, \ldots, s_k^\rho\}$ represents any term formed by combining all the terms $s_0$, $\ldots$, $s_k$ (possibly with repetitions) via the non-terminal $\mathsf{D}_\rho$. If $S$ is a finite set of terms of the same type, it follows that $S \to^* s$ for each $s \in S$.

Notice that there are no start symbols in $\mathscr{H}$. In this regard we may consider the individual Herbrand scheme $\mathscr{H}_\pi$ as obtained from $\mathscr{H}$ by specifying an appropriate set of start symbols. The new 'hat' non-terminals do not play a role in viewing $\mathscr{H}$ as a universal Herbrand scheme. Rather, they become useful in 'transferring' non-terminals lacking their final argument through applications of permutation. For example, the following partial proof (where we assume $u(A) > 0$) gives rise to the production rules on the right:

$$
\begin{array}{ll}
\dfrac{\mathsf{p} \; \dfrac{\pi_1 \vdash A(r/v), B}{\pi_0 \vdash B, A(r/v)}}{\exists_r \; \dfrac{}{\pi \vdash B, \exists v A}}
&
\begin{array}{l}
\mathsf{N}_{\pi_0}^i a x z \to \mathsf{N}_{\pi_1}^{1-i} a z x \\[4pt]
\mathsf{N}_\pi^1 a x z \to \langle r \cdot a, \mathsf{N}_{\pi_0}^1 x \rangle \\[4pt]
\mathsf{N}_\pi^0 a x z \to \mathsf{N}_{\pi_0}^0 a x \big( z(\mathsf{N}_\pi^1 a x z) \big)
\end{array}
\end{array}
$$

yielding the derivation $\mathsf{N}_\pi^0 a x z \to^* \mathsf{N}_{\pi_1}^1 a \big( z \langle r \cdot a, \mathsf{N}_{\pi_0}^1 a x \rangle \big) x$. The derivation cannot be extended as it stands because $\mathsf{N}_{\pi_0}^1$ lacks an argument, meaning that it is not formally possible to express the term $\mathsf{N}_\pi^0 a x z$ by reference to the proof $\pi_1$ only without instantiating $x$ and $z$ by concrete terms. However, $\mathsf{N}_{\pi_0}^1 a x$ is extensionally equal to the term $\hat{\mathsf{N}}_{\pi_1}^0 a x$, allowing us to equate $\mathsf{N}_\pi^0 a x z$ with the term $\mathsf{N}_{\pi_1}^1 a \big( z \langle r \cdot a, \hat{\mathsf{N}}_{\pi_1}^0 a x \rangle \big) x$ for any choice of $a$, $x$ and $z$. Equations such as these are useful in the close examination of the cut elimination process carried out in the sections below.

In the previous section a natural subsumption and equivalence relation on terms was introduced given by equating terms that induce the same language in all contexts. In the context of the universal Herbrand scheme $\mathscr{H}$, this subsumption is given by $r \prec s$ which holds just if $r, s : \rho$ for some $\rho$ and, for every $\mathscr{H} \cup \{x^\rho\}$-term $t$ of basic type, whenever $t(r/x) \to^* r_0$ for a $\Sigma$-term $r_0$, then $t(s/x) \to^* s_0$ for some $\Sigma$-term $s_0$ such that $r_0^\circ = s_0^\circ$. The corresponding equivalence relation $\asymp$ is defined by $r \asymp s$ iff $r \prec s \prec r$. The following properties of the relations $\prec$ and $\asymp$ were remarked in the last section.

**Lemma 5.2.** *Let $r, s : \rho$ be $\mathscr{H}$-terms of the same type and $S$ a finite set of terms of pair type $\sigma = \sigma_0 \times \cdots \times \sigma_l$.*

1. *If $r \to s$ then $s \prec r$. If, in addition, the reduction follows from a production rule for a deterministic non-terminal then $r \asymp s$ .*

2. *If $r$ and $s$ are representations of the same finite set of $\mathscr{H}$-terms then $r \asymp s$.*

*3. If $r = \mathsf{F}s_0 \cdots s_{k-1}x^\sigma$ for a non-terminal $\mathsf{F}$ with production rule*

$$\mathsf{F}x_0 \cdots x_{k-1}\langle x_k, \ldots, x_{k+l}\rangle \to t$$

*then*

$$r(S/x) \asymp \{r(s/x) \mid s \in S\} \asymp \{t((s_0 \ldots, s_{k+l})/(x_0, \ldots, x_{k+l})) \mid \langle s_k, \ldots, s_{k+l}\rangle \in S\}.$$

*4. If the co-domain of $\rho$ is $\epsilon$ then $r \asymp c_\rho$.*

*Proof.* Properties 1–3 are straight-forward, though for 3 we note that only deterministic non-terminals have production rules that invoke pattern-matching. 4 generalises Lemma 3.29 and is proved by induction on $\rho$ and $r$, noting that $\langle\rangle \cdot a \asymp \langle\rangle$ for any substitution $a$. $\qquad\square$

Both relations can be extended to proofs by setting $\pi' \prec \pi$ if $\pi$ and $\pi'$ have the same end-sequent and $\mathsf{N}^i_{\pi'} \prec \mathsf{N}^i_\pi$ for each $i$. For many cases of $\pi \rightsquigarrow \pi'$ indeed $\pi' \prec \pi$ (or even $\pi' \asymp \pi$), from which we may immediately deduce $L(\pi') \subseteq L(\pi)$ (resp. $L(\pi') = L(\pi)$). However, there exist reductions $\pi \rightsquigarrow \pi'$ for which $L(\pi') \subseteq L(\pi)$ but $\pi' \not\prec \pi$. These scenarios all concern reductions that interact with quantifiers and alter the contexts in which explicit substitutions appear in derivations. To establish language preservation also in these cases we introduce a coarser relation $\prec$, denoted $\sqsubset$, based on quantifying over contexts of a particular syntactic shape that is preserved by $\mathscr{H}$-derivations. Such terms, which we call *normal $\mathscr{H}$-terms*, are defined below.

## 5.1 Normal Terms and Subsumption

In order to focus on the impact of substitutions in $\mathscr{H}$-terms it is necessary to introduce a notion of free and bound occurrences of $\Sigma$-symbols in these terms where, recall, $\Sigma$ is the signature of first-order logic. The *free* symbols of a basic $\Sigma$-term are simply the $\Sigma$-symbols that occur in the term; $\Sigma$-terms have no *bound* symbols. For a basic $\mathscr{H}$-term $t$, the *free* symbols of $t$ are the $\Sigma$-symbols occurring in $t$ combined with the $\Sigma$-symbols occurring in any proof $\pi$ for which a non-terminal $\mathsf{N}^i_\pi$ or $\hat{\mathsf{N}}^i_\pi$ appears in $t$; the *bound* symbols of $t$ are the eigenvariables of the proofs which occur leftmost in $t$. For non-basic terms, substitutions and substitution stacks are interpreted as contributing to the set of bound symbols, and limiting the set of free symbols in the natural way. Explicitly, for a substitution stack $a : \varsigma$ and $\mathscr{H}$-term $r : \rho$ we define

$$Bd(r) = \begin{cases} \emptyset, & \text{if } r \text{ is a } \Sigma\text{-term or } r = \mathsf{c}_\rho, \\ EV(\pi), & \text{if } r = \mathsf{N}^i_\pi \text{ or } r = \hat{\mathsf{N}}^i_\pi \text{ for some } i, \\ Bd(s) \cup Bd(t), & \text{if } r = \langle s, t\rangle, \\ Bd(s) \cup Bd(t), & \text{if } r = \mathsf{D}_\rho st, \\ Bd(s) \cup Bd(a), & \text{if } r = sa \text{ or } r = s \cdot a \text{ for } a : \varsigma, \\ Bd(s), & \text{if } r = st \text{ and } t : \tau \text{ where } \tau \neq \varsigma, \end{cases}$$

$$Bd(a) = \begin{cases} \emptyset, & \text{if } a = \bot, \\ \{\alpha\} \cup Bd(b), & \text{if } a = [\alpha \hookleftarrow s]b, \end{cases}$$

$$Fr(r) = \begin{cases} \emptyset, & \text{if } r = \mathsf{c}_\rho \text{ or } r = \mathsf{D}_{\rho'}, \\ \{r\}, & \text{if } r \in \Sigma, \\ Fr(\pi), & \text{if } r = \mathsf{N}_\pi^i \text{ or } r = \hat{\mathsf{N}}_\pi^i \text{ for some } i, \\ Fr(s) \cup Fr(t), & \text{if } r = \langle s, t \rangle, \\ (Fr(s) \setminus Bd(a)) \cup Fr(a), & \text{if } r = sa \text{ or } r = s \cdot a \text{ for } a : \varsigma, \\ Fr(s) \cup Fr(t), & \text{if } r = st \text{ and } t : \tau \text{ where } \tau \neq \varsigma, \end{cases}$$

$$Fr(a) = \begin{cases} \emptyset, & \text{if } a = \bot, \\ (Fr(s) \setminus Bd(b)) \cup Fr(b), & \text{if } a = [\alpha \hookleftarrow s]b, \end{cases}$$

$EV(\pi)$ denotes the set of eigenvariables in the proof $\pi$, and $Fr(\pi)$ the set of all non-eigenvariable $\Sigma$-symbols occurring in the $\pi$. Notice that $Bd(\mathsf{N}_\pi^i)$ and $Fr(\mathsf{N}_\pi^j)$ are disjoint sets by definition.

**Definition 5.3** (Normal terms)**.** A *normal term* is an $\mathscr{H}$-term $r$ satisfying:

1. if $a$ is substitution stack which is a sub-term of $r$ then $Bd(a) \cap Fr(a) = \emptyset$,

2. if $st$ is an application which is a sub-term of $r$ then $Bd(s) \cap Fr(t) = \emptyset$,

3. if $s \cdot a$ is a substitution occurring as a sub-term of $r$ then $s$ is of simple type.

As mentioned at the beginning of this section, the aim of the above definition is to provide a class of terms for which we can examine a more refined subsumption relation on $\mathscr{H}$-terms that captures both language inclusion and equality for a wide range of cut reduction rules. The subsumption relation that achieves this is essentially the restriction of $\prec$ that only quantifies over normal contexts.

**Definition 5.4** (Subsumption)**.** Given normal $\mathscr{H}$-terms $r, s : \rho$ of the same type, $s$ *subsumes* $r$, in symbols $r \sqsubseteq s$, just if, for every $\mathscr{H} \cup \{x^\rho\}$-term $t$ of basic type such that $t(r/x)$ and $t(s/x)$ are both normal, whenever $t(r/x) \to^* u$ for a $\Sigma$-term $u$ then $t(s/x) \to^* v$ for some $\Sigma$-term $v$ satisfying $u^\circ = v^\circ$. Define $r \sim s$ if $r \sqsubseteq s$ and $s \sqsubseteq r$.

Clearly, for normal terms $r$ and $s$, $r \prec s$ implies $r \sqsubseteq s$, and $r \asymp s$ implies $r \sim s$. Hence, if $\pi, \pi'$ are two regular proofs of a $\Sigma_1$ sequent $\Gamma$ and $\mathsf{S}_{\pi',i} \sqsubseteq \mathsf{S}_{\pi,i}$ for every $i < |\Gamma|$ then $L(\pi') \subseteq L(\pi)$. However, what we require is the more general property that if for every $\mathscr{H}$-derivation $\mathsf{S}_{\pi',i} \to^* u$ of a $\Sigma$-term there exists $\mathscr{H}$-terms $r$, $s$ and $t$ such that $\mathsf{S}_{\pi',i} \to^* t(r/x) \to^* u$, $\mathsf{S}_{\pi,i} \to^* t(s/x)$ and $r \sqsubseteq s$, then we may conclude $L(\pi') \subseteq L(\pi)$. This result holds trivially for $\prec$ in place of $\sqsubseteq$. For it to work for subsumption, the terms $r$, $s$ and $t$ must all be normal, i.e., we require

**Lemma 5.5.** *If* $\mathsf{S}_{\pi,i} \to^* s$ *then* $s$ *is a normal term.*

The proof of Lemma 5.5 is not difficult but requires some technical observations concerning the preservation of free and bound symbols through $\mathscr{H}$-derivations, so will be postponed to the next section (Lemma 5.10). Another key lemma concerning derivations is the following which, when combined with Lemma 5.5, implies that all derivations from start symbols lead to $\Sigma$-terms.

**Lemma 5.6.** *If $r$ is a normal term of prime type and not a $\Sigma$-term then $r \to s$ for some term $s$.*

There are two scenarios in which a term $r$ that is not a $\Sigma$-term cannot be reduced. First, the reduction of $r$ may require pattern-matching on an argument that cannot be reduced to a pair, such as $\langle s, t \rangle \cdot a$. Second, $r$ may contain a sub-term $\mathsf{F} r_1 \cdots r_k \cdot a$ for which the arity of $\mathsf{F}$ is greater than $k$. Normality prevents either scenario from occurring (Lemma 5.11). As a consequence, terms of pair type always reduce to pairs (Lemma 5.12) and pattern-matching will not block derivations. Thus,

**Lemma 5.7** (Finite basis lemma). *For every normal $\mathscr{H}$-term $r$ of pair type there exists terms $\langle s_0, t_0 \rangle, \ldots, \langle s_k, t_k \rangle$ such that $r \sim \{\langle s_i, t_i \rangle \mid i \le l\}$.*

We now turn to proving these three lemmas, and establishing some further properties of subsumption. The reader may wish to proceed directly to Section 5.3 at this point and refer back to Section 5.2 as needed.

## 5.2 Technical Lemmas

In order to establish Lemma 5.5, we require two observations on free and bound symbols in normal terms. The effect is to reduce the work to checking that production rules locally preserve normality (in contrast to within arbitrary contexts).

**Lemma 5.8.** *If $r(s/x)$ is a normal term, $t$ is a normal term of the same type as $s$, $Fr(t) \subseteq Fr(s)$ and $Bd(t) \subseteq Bd(s)$ then $r(t/x)$ is normal.*

*Proof.* By definition. $\qquad\square$

**Lemma 5.9.** *If $\mathsf{F} x_0 \cdots x_{k-1} \langle x_k, \ldots, x_{k+l} \rangle \to t$ is a production rule of $\mathscr{H}$, and $r_0, \ldots, r_{k+l}$ are such that $s = \mathsf{F} r_0 \cdots r_{k-1} \langle r_k, \ldots, r_{k+l} \rangle$ is normal, then $Fr(t(\vec{r}/\vec{x})) \subseteq Fr(s)$ and $Bd(t(\vec{r}/\vec{x})) \subseteq Bd(s)$.*

*Proof.* We examine two particular cases, namely the quantifier rules, and leave the remaining for the reader to check. Consider an instance of the production rule for $\forall_\alpha$ for a single eigenvariable:

$$\mathsf{N}^i_\pi a r_0 \cdots r_{k-1} \langle s, r_k \rangle \to \mathsf{N}^i_{\pi_0}([\alpha \leftarrow s]a) r_0 \cdots r_k$$

where $r_0, \ldots, r_k$ and $s$ are terms of suitable type and number, and $a : \varsigma$ is a substitution stack. Let $m$ and $n$ abbreviate the left- and righthand term in the above rule respectively. Assume $m$ is a normal term.

$$\begin{aligned}
Bd(m) &= Bd(\mathsf{N}^i_\pi a) \\
&= EV(\pi) \cup Bd(a) \\
&= EV(\pi_0) \cup \{\alpha\} \cup Bd(a) \\
&= Bd(n)
\end{aligned}$$

Concerning free symbols, we have

$$\begin{aligned}
Fr(m) &= (Fr(\pi) \setminus Bd(a)) \cup Fr(a) \cup Fr(r_0, \ldots, r_k, s) \\
Fr(n) &= (Fr(\pi_0) \setminus Bd([\alpha \leftarrow s]a)) \cup Fr([\alpha \leftarrow s]a) \cup Fr(r_0, \ldots, r_k) \\
&= (Fr(\pi_0) \setminus (\{\alpha\} \cup Bd(a))) \cup (Fr(s) \setminus Bd(a)) \cup Fr(a) \cup Fr(r_0, \ldots, r_k)
\end{aligned}$$

where $Fr(u_0, \ldots, u_l) = \bigcup_{i \leq l} Fr(u_i)$. By normality of $m$, $Fr(s) \setminus Bd(a) = Fr(s)$ and as $Fr(\pi) = Fr(\pi_0) \setminus \{\alpha\}$, so $Fr(n) \subseteq Fr(m)$.

For production rules resulting from the inference rule $\exists_s$, suppose

$$\mathsf{N}_\pi^i a r_0 \cdots r_k \to \mathsf{N}_{\pi_0}^i a r_0 \cdots r_{k-1} (r_k \langle s \cdot a, \mathsf{N}_{\pi_0}^k a r_0 \cdots r_{k-1} \rangle) \tag{13}$$

for suitable terms $r_0$, $\ldots$, $r_k$ and $a$. Recall that $s$ is the $\Sigma$-term instantiating the existential quantifier in the active formula of $\pi_0$. By our regularity condition on proofs, $Fr(s) \subseteq Fr(\pi)$, so, letting $m$ and $n$ denote the left and right side of the reduction in (13), we have

$$\begin{aligned}
Fr(n) &= Fr(\mathsf{N}_{\pi_0}^i a) \cup Fr(r_0, \ldots, r_k) \cup Fr(s \cdot a) \cup Fr(\mathsf{N}_{\pi_0}^k a) \\
&\subseteq Fr(\mathsf{N}_\pi^i a) \cup Fr(r_0, \ldots, r_k) \\
&= Fr(m). \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \square
\end{aligned}$$

The following lemma implies Lemma 5.5.

**Lemma 5.10.** *If $r \to s$ and $r$ is normal then $s$ is normal.*

*Proof.* By the previous two lemmas, it suffices to show that every production rule of $\mathscr{H}$ locally preserves normality. As in the previous proof, we offer the argument for some particular cases of the two quantifier rules. Let the derivation

$$\mathsf{N}_\pi^i a r_1 \cdots r_k \langle s, r_{k+1} \rangle \to \mathsf{N}_{\pi_0}^i ([\alpha \leftarrowtail s] a) r_1 \cdots r_k r_{k+1}$$

arise from an inference $\forall_\alpha$, with $m$ and $n$ denoting the left and righthand term. Assume $m$ is normal. In particular,

$$(EV(\pi) \cup Bd(a)) \cap (Fr(a) \cup Fr(s)) = \emptyset \tag{14}$$

We first show that for every application $s't'$ occurring in $n$, $Bd(s') \cap Fr(t') = \emptyset$. This is evident if $s't'$ is a sub-term of $a$, $r_1$, $\ldots$, $r_k$, $s$ or $t$. Moreover, it holds for the case $s' = \mathsf{N}_{\pi_0}^i$ and $t' = [\alpha \leftarrowtail s] a$ because

$$\begin{aligned}
Bd(s') \cap Fr(t') &= EV(\pi_0) \cap ((Fr(s) \setminus Bd(a)) \cup Fr(a)) \\
&\subseteq EV(\pi) \cap (Fr(s) \cup Fr(a)) \\
&= \emptyset
\end{aligned}$$

and for $s' = \mathsf{N}_{\pi_0}^i ([\alpha \leftarrowtail s] a) r_1 \cdots r_j$, $t' = r_{j+1}$ $(j \leq k)$ because $Bd(\mathsf{N}_{\pi_0}^i ([\alpha \leftarrowtail s] a)) \subseteq Bd(\mathsf{N}_\pi^i a)$. The other requirement to check for normality is that the sets $Bd([\alpha \leftarrowtail s] a)$ and $Fr([\alpha \leftarrowtail s] a)$ are disjoint, but this follows from (14), given that $Fr([\alpha \leftarrowtail s] a) \subseteq Fr(s) \cup Fr(a)$ and $Bd([\alpha \leftarrowtail s] a) \subseteq EV(\pi) \cup Bd(a)$. Hence $n$ is normal.

The second case we consider is the one-step derivation

$$\mathsf{N}_\pi^i a r_0 \cdots r_k \to \mathsf{N}_{\pi_0}^i a r_0 \cdots r_{k-1} (r_k \langle s \cdot a, \mathsf{N}_{\pi_0}^k a r_0 \cdots r_{k-1} \rangle)$$

due to the inference $\exists_s$. Suppose $\mathsf{N}_\pi^i a r_0 \cdots r_k$ is a normal term so, in particular, $Bd(\mathsf{N}_\pi^i a)$ is disjoint from $Fr(r_i)$ for each $i \leq k$. In this case it only needs establishing that

$$Bd(\mathsf{N}_{\pi_0}^i a r_0 \cdots r_{k-1}) \cap Fr(r_k \langle s \cdot a, \mathsf{N}_{\pi_0}^k a r_0 \cdots r_{k-1} \rangle) = \emptyset$$

i.e., that

$$Bd(\mathsf{N}_{\pi_0}^i a) \cap \left( \bigcup_{i \le k} Fr(r_i) \cup Fr(s \cdot a) \cup Fr(\mathsf{N}_{\pi_0}^k a) \right) = \emptyset,$$

as all other cases follow immediately from normality of $\mathsf{N}_\pi^i a r_0 \cdots r_k$. But by regularity of $\pi$, $Fr(s) \subseteq Fr(\pi)$, so we have

$$Bd(\mathsf{N}_{\pi_0}^i a) = EV(\pi) \cup Bd(a)$$
$$Fr(s \cdot a) \subseteq (Fr(\pi) \setminus Bd(a)) \cup Fr(a)$$
$$Fr(\mathsf{N}_{\pi_0}^k a) = (Fr(\pi) \setminus Bd(a)) \cup Fr(a)$$

and, as $Bd(a)$ is disjoint from $Fr(a)$ and $EV(\pi)$ is disjoint from $Fr(\pi) \cup Fr(a)$, we are done. $\qquad\square$

We now turn to the task of proving Lemma 5.6 which, as explained earlier, follows from the next two lemmas. The first lemma characterises the form normal $\mathscr{H}$-terms may take, and will be useful in the subsequent analysis of derivations in Herbrand schemes.

**Lemma 5.11.** *If $r : \rho$ is a normal $\mathscr{H}$-term and $\rho$ is a basic type whose co-domain is a pair $\sigma \times \tau$ in which $\sigma$ is simple and $\tau$ is not simple, then either $r = \langle s, t \rangle$ for a $\Sigma$-term $s$ and $\mathscr{H}$-term $t$, or $r = \mathsf{F}r_1 \cdots r_k$ for some non-terminal $\mathsf{F}$ and terms $r_1, \ldots, r_k$.*

*Proof.* By induction on $r$. Let $r : \rho = \rho_1 \to \cdots \to \rho_l \to \sigma \times \tau$ be a normal $\mathscr{H}$-term satisfying the hypothesis of the lemma. Since $\rho$ is not a simple type, $r$ is not a $\Sigma$-symbol nor of the form $s \cdot a$ for a substitution stack $a$ (by definition of normal terms). This leaves three cases: i) $l = 0$ and $r = \langle s, t \rangle$ for $s : \sigma$ and $t : \tau$; ii) $r = st$ for $s : \tau' \to \rho$ and $t : \tau'$; or iii) $r$ is a non-terminal of $\mathscr{H}$. If (i), as $\sigma$ is simple, $s$ is a $\Sigma$-term by Lemma 3.27 and we are done. In case (ii), suppose $r = st$ is an application and $s : \sigma' = \tau' \to \rho$ and $t : \tau'$. If $\sigma'$ is not basic then Lemma 3.28 implies $s = \mathsf{F}$ for some $\pi$ and $i$, whence $r = \mathsf{F}t$. On the other hand, if $\sigma'$ is basic the induction hypothesis applies and $s = \mathsf{F}r_1 \cdots r_k$ for terms $r_1, \ldots, r_k$, and so similarly for $r$. So we are done. $\qquad\square$

**Lemma 5.12.** *If $r : \iota \times \rho$ is a normal $\mathscr{H}$-term of pair type but not a pair then $r \to s$ for some $\mathscr{H}$-term $s$.*

*Proof.* Assume to the contrary that $r : \iota \times \rho$ is an $\mathscr{H}$-term which is not a pair and that there is no $s$ such that $r \to s$. Without loss of generality assume $r$ is minimal in length. By Lemma 5.11, $r = \mathsf{F}r_1, \ldots, r_k$ for some non-terminal $\mathsf{F}$ and terms $r_1, \ldots, r_k$. It follows that $\mathsf{F} \ne \mathsf{c}_\sigma$ for any $\sigma$ as otherwise $r \to \langle \mathsf{c}, \mathsf{c}_\rho \rangle$. Also $\mathsf{F} \ne \mathsf{D}_{\iota \times \rho}$ (as then $r \to r_1$) and $\mathsf{F} \ne \hat{\mathsf{N}}_\pi^i$ for any $\pi$ and $i$. So $\mathsf{F} = \mathsf{N}_\pi^i$ for some $\pi$ and $i$. The fact that $r$ is not reducible means that the production rule for $\mathsf{N}_\pi^i$ requires pattern-matching on the final argument. But then $r_k : \iota \times \sigma$ for some $\sigma$, is not a pair and is not reducible, contradicting minimality of $r$. $\qquad\square$

In the remainder of this section we present some basic properties of the subsumption relation which will be needed in the next section, starting with a proof of the Finite Basis Lemma.

*Proof of Lemma 5.7.* Let $r : \rho$ be a $\mathscr{H}$-term where $\rho = \sigma \times \tau$. Without loss of generality, we may assume $r \neq \mathsf{D}_\rho r_0 r_1$ for any $r_0$ and $r_1$. If $r$ has the form $\langle s, t \rangle$ then trivially $r \sim \{\langle s, t \rangle\}$ and if $r = \mathsf{c}_\rho$ then $r \sim \{\langle \mathsf{c}_\sigma, \mathsf{c}_\tau \rangle\}$. Otherwise, Lemma 5.11 implies that $r \sim \mathsf{N}_\pi^i a r_0 \cdots r_k$ for some $\pi$, $i$, $a$, $r_0$, ..., $r_k$. An induction on $\pi$ determines terms $s_0, \ldots, s_l$ and $t_0, \ldots, t_l$ such that $r \sim \{\langle s_j, t_j \rangle \mid j \leq l\}$. Note that there is no issue with pattern-matching stopping derivations from fully writing out, as Lemma 5.12 demonstrates. $\qquad\square$

**Lemma 5.13.** *If $t(r/x)$ and $t(s/x)$ are normal terms and $r \sqsubset s$ then $t(r/x) \sqsubset t(s/x)$.*

*Proof.* Direct consequence of the definition. $\qquad\square$

**Lemma 5.14.** *If $ru \sqsubset su$ for every term $u$ then $r \sqsubset s$.*

*Proof.* For every derivation $t(r/x) \to^* r_0$ of a $\Sigma$-term, there are terms $t'$, $u_0$, ..., $u_k$ such that $t(r/x) \to^* t'((ru_i)_{i \leq k}/\vec{x}) \to^* r_0$ and $t(s/x) \to^* t'((su_i)_{i \leq k}/\vec{x})$. Since normality is preserved through derivations, we are done. $\qquad\square$

**Lemma 5.15.** *Let $r$ be a basic $\mathscr{H}$-term and $a$ be a substitution stack over $\Sigma$. Then,*

   *i)* $Fr(r(s/\alpha) \cdot a) \subseteq Fr(r \cdot ([\alpha \hookleftarrow s]a))$,

   *ii)* $Bd(r(s/\alpha) \cdot a) \subseteq Bd(r \cdot ([\alpha \hookleftarrow s]a))$ *provided $s$ is basic,*

   *iii)* *if $Fr(r) \cap Bd(a) = \emptyset$ then $r^a = r$,*

   *iv)* *if $\alpha \notin Bd(a) \cup Fr(a)$ then $r(s/\alpha)^a = r^a(s^a/\alpha)$ provided $s$ is basic.*

*Proof.* By induction on $r$ and $a$. $\qquad\square$

**Lemma 5.16.** *Let $r : \rho$ and $a : \varsigma$ be $\Sigma$-terms, $s : \sigma$ a basic $\Sigma$-term, $\alpha^\rho \in \Sigma$ and $\pi$ a regular proof. Under the assumption that $\mathsf{N}_\pi^i([\alpha \hookleftarrow r]a)$ is normal the following hold.*

   *1.* $s \cdot ([\alpha \hookleftarrow r]a) \sim s(r/\alpha) \cdot a$,

   *2.* $r \sim r^\circ$,

   *3.* $[\alpha \hookleftarrow r]a \sim [\alpha \hookleftarrow r \cdot a]a$,

   *4.* *If $EV(\pi) \cap Fr(r) = \emptyset$ then $\mathsf{N}_\pi^i([\alpha \hookleftarrow r]a) \sim \mathsf{N}_{\pi(r^\circ/\alpha)}^i a$,*

   *5.* *If $\alpha \notin Fr(\pi)$ then $\mathsf{N}_\pi^i a \sim \mathsf{N}_\pi^i([\alpha \hookleftarrow r]a)$.*

*Proof.* 1 is proved via induction on the basic term $s$. That $r$ and $a$ are $\Sigma$-terms is necessary for showing $s(r/\alpha) \cdot a \sqsubset s \cdot ([\alpha \hookleftarrow r]a)$. 2 follows from 1 by induction on $r$. Regarding 3, Lemma 5.15(i,ii) imply $Fr(a^\circ) \subseteq Fr(a)$ and $Bd(a^\circ) \subseteq Bd(a)$, so $\alpha \notin Fr(a^\circ)$ by normality. Hence, if $t$ is a basic term then

$$
\begin{aligned}
t \cdot [\alpha \hookleftarrow r]a &\sim t(r^\circ/\alpha)^{a^\circ} \sim t^{a^\circ}((r^\circ)^{a^\circ}/\alpha) \\
&\sim t^{a^\circ}(((r^\circ)^{a^\circ})^{a^\circ}/\alpha) \\
&\sim t((r^\circ)^{a^\circ}/\alpha)^{a^\circ} \\
&\sim t \cdot [\alpha \hookleftarrow r \cdot a]a.
\end{aligned}
$$

The first and last equivalence are applications of 2; the second and fourth equivalence are consequences of Lemma 5.15(iv); and the third equivalence uses Lemma 5.15(iii) and the fact that $Fr(a^\circ) \cap Bd(a^\circ) = \emptyset$. Via 2 the above holds for $t$ an arbitrary $\Sigma$-term, and from there generalises to deduce $[\alpha \leftarrowtail r]a \sim [\alpha \leftarrowtail r \cdot a]a$.

4 is derived by induction on $\pi$. By 2 we may assume $r$ is a basic term, i.e., $r = r^\circ$. In the base case, where $\pi$ is an axiom, the equivalence is trivial as $\mathsf{N}_\pi^i([\alpha \leftarrowtail r]a)r_1 r_2 \sim \mathsf{N}_{\pi^{(r/\alpha)}}^i a r_1 r_2$ for any choice of $r_1$ and $r_2$ of appropriate type. The induction step is straightforward except in the case of quantifier rules. If $\pi$ ends in the inference

$$\exists_{\vec{s}} \; \frac{\pi_0 \vdash \Gamma, A(\vec{s}/\vec{v})}{\pi \vdash \Gamma, \exists \vec{v} A}$$

where $\vec{s} = (s_j)_{j \leq k}$ then we have, if $i = |\Gamma|$, $b = [\alpha \leftarrowtail r]a$ and $r_1$, ..., $r_{|\Gamma|}$ and $t$ are suitable normal terms,

$$\mathsf{N}_\pi^{|\Gamma|} b r_1 \cdots r_{|\Gamma|} t \sim \langle s_0 \cdot b, \ldots, s_k \cdot b, \mathsf{N}_{\pi_0}^{|\Gamma|} b r_1 \cdots r_{|\Gamma|} \rangle$$
$$\sim \langle s_0(r/\alpha) \cdot a, \ldots, s_k(r/\alpha) \cdot a, \mathsf{N}_{\pi_0^{(r/\alpha)}}^{|\Gamma|} a r_1 \cdots r_{|\Gamma|} \rangle$$
$$\sim \mathsf{N}_{\pi^{(r/\alpha)}}^{|\Gamma|} a r_1 \cdots r_{|\Gamma|} t$$

where the second equivalence due to the induction hypothesis for $\mathsf{N}_{\pi_0}^{|\Gamma|}([\alpha \leftarrowtail r]a)$. The case $i < |\Gamma|$ is similar. For applications of the $\forall$ inferences, we consider the inference

$$\forall_{\vec{\beta}} \; \frac{\pi_0 \vdash \Gamma, A(\vec{\beta}/\vec{v})}{\pi \vdash \Gamma, \forall \vec{v} A}$$

for $\vec{\beta} = (\beta_j^\iota)_{j \leq k}$. By normality of $\mathsf{N}_\pi^i([\alpha \leftarrowtail r]a)$ and the assumption that $Fr(r) \cap EV(\pi) = \emptyset$, it follows that $\beta_j \notin Fr(a) \cup Fr(r)$ for each $j$. Let $r_1$, ..., $r_{|\Gamma|}$, $s_0$, ..., $s_k$ and $t$ be such that the term $n := \mathsf{N}_\pi^i([\alpha \leftarrowtail r]a)r_1 \cdots r_{|\Gamma|} \langle s_0, \ldots, s_k, t \rangle$ is well-typed and normal. In particular, $\alpha, \beta_0, \ldots, \beta_k \notin Fr(\langle s_0, \ldots, s_k, t \rangle)$. Moreover, as $s_j$ has simple type for each $j \leq k$, Lemma 3.27 implies that $\vec{s}$ is a sequence of $\Sigma$-terms. Then assuming $\alpha \notin \{\beta_j \mid j \leq k\}$, and writing $[\vec{\beta} \leftarrowtail \vec{s}]b$ in place of $[\beta_0 \leftarrowtail s_0] \cdots [\beta_k \leftarrowtail s_k]b$, we have

$$n \sim \mathsf{N}_{\pi_0}^i([\vec{\beta} \leftarrowtail \vec{s}][\alpha \leftarrowtail r]a)r_1 \cdots r_{|\Gamma|} t$$
$$\sim \mathsf{N}_{\pi_0^{(\vec{s}/\vec{\beta})(r/\alpha)}}^i a r_1 \cdots r_{|\Gamma|} t$$
$$\sim \mathsf{N}_{\pi_0^{(r/\alpha)(\vec{s}/\vec{\beta})}}^i a r_1 \cdots r_{|\Gamma|} t$$
$$\sim \mathsf{N}_{\pi_0^{(r/\alpha)}}^i([\vec{\beta} \leftarrowtail \vec{s}]a)r_1 \cdots r_{|\Gamma|} t$$
$$\sim \mathsf{N}_{\pi^{(r/\alpha)}}^i a r_1 \cdots r_{|\Gamma|} \langle s_0, \ldots, s_k, t \rangle.$$

The third equivalence holds since $\alpha \notin Fr(\langle s_0, \ldots, s_k, t \rangle)$ and $\beta_j \notin Fr(r)$ for any $j$. If

$\alpha = \beta_j$ then $\pi^{(r/\alpha)} = \pi$ and, using again that $\alpha \notin Fr(\langle s_0, \ldots, s_k \rangle)$, we have

$$
\begin{aligned}
n &\sim \mathsf{N}^i_{\pi_0}([\vec{\beta} \leftarrowtail \vec{s}][\alpha \leftarrowtail r]a)r_1 \cdots r_{|\Gamma|}t \\
&\sim \mathsf{N}^i_{\pi_0^{(\vec{s}/\vec{\beta})(r/\alpha)}}ar_1 \cdots r_{|\Gamma|}t \\
&\sim \mathsf{N}^i_{\pi_0^{(\vec{s}/\vec{\beta})}}ar_1 \cdots r_{|\Gamma|}t \\
&\sim \mathsf{N}^i_{\pi_0}([\vec{\beta} \leftarrowtail \vec{s}]a)r_1 \cdots r_{|\Gamma|}t \\
&\sim \mathsf{N}^i_{\pi}ar_1 \cdots r_{|\Gamma|}\langle s_0, \ldots, s_k, t\rangle \\
&\sim \mathsf{N}^i_{\pi^{(r/\alpha)}}ar_1 \cdots r_{|\Gamma|}\langle s_0, \ldots, s_k, t\rangle.
\end{aligned}
$$

5 is a special case of 4. $\qquad \square$

Let $\pi$ be a proof with end-sequent $A_1, \ldots, A_m, B, C_1, \ldots, C_n$ and suppose $a : \varsigma$, $r_i : \tau^*_{A_i}$, $s : \tau^*_B$ and $t_j : \tau^*_{C_j}$ are terms for each $1 \le i \le m$ and $1 \le j \le n$ such that $\mathsf{N}^m_\pi a\vec{r}s\vec{t}$ is normal. Let $\rho = \tau^*_B$.

**Lemma 5.17.** *If $B$ is prenex $\Sigma_1$ then*

$$
\mathsf{N}^m_\pi a\vec{r}s\vec{t} \sim \mathsf{N}^m_\pi a\vec{r}\mathsf{c}_\rho \vec{t}.
$$

*Proof.* Since $B$ is $\Sigma_1$, $\tau^*_B = \tau_B \to \epsilon$. Lemma 5.2 completes the proof. $\qquad \square$

**Lemma 5.18.** *If $e(B) > 0$ and there are no applications of contraction to $B$ in $\pi$ then*

$$
\mathsf{N}^m_\pi a\vec{r}s\vec{t} \sim \mathsf{N}^m_\pi a\vec{r}\mathsf{c}_\rho \vec{t}.
$$

**Lemma 5.19.** *If the final inference in $\pi$ is an application of $\mathsf{p}$ with sub-proof $\pi' \vdash A_1, \ldots, A_m, C_1, B, C_2, \ldots, C_n$ then for each $j \in [0, m) \cup [m + 2, m + n]$,*

$$
\mathsf{N}^j_\pi ar_1 \cdots r_m st_1 \cdots t_n \sim \mathsf{N}^j_{\pi'}ar_1 \cdots r_m t_1 st_2 \cdots t_n
$$

The proofs of the final two lemmas proceed by induction on $\pi$.

## 5.3 Language Preservation for Gentzen-style Cut Elimination

Let $\pi$ and $\pi'$ be regular proofs of some sequent $\Gamma$. We say that $\pi$ *subsumes* $\pi'$, in symbols $\pi' \sqsubset \pi$, if $\mathsf{N}^i_{\pi'} \sqsubset \mathsf{N}^i_\pi$ for every $i < |\Gamma|$. If $\pi$ and $\pi'$ each subsumes the other then $\pi$ and $\pi'$ are *equivalent*, in symbols $\pi \sim \pi'$. As an immediate consequence of the definition we have

**Lemma 5.20.** *Suppose $\pi$ and $\pi'$ are proofs of the same $\Sigma_1$ sequent. If $\pi' \sqsubset \pi$ then $L(\pi') \subseteq L(\pi)$.*

Herbrand schemes have the property that their languages are invariant under many basic proof transformations. The first example we give concerns the operation of substitution in proofs:

**Lemma 5.21.** *Suppose $\pi$ and $\pi'$ are proofs with the same end-sequent such that $\pi'$ is the result of replacing a sub-proof $\pi_0$ of $\pi$ by $\pi'_0$. If $\pi_0 \sqsubset \pi'_0$ then $\pi \sqsubset \pi'$.*

*Proof.* Let $\pi$, $\pi_0$, $\pi'$ and $\pi_0'$ be as in the statement. We assume $\pi$ and $\pi'$ have the same end-sequent, say $\Gamma$. Given a subproof $\hat\pi$ of $\pi$ which is not a proper subproof of $\pi_0$, let $\hat\pi'$ denote the corresponding subproof of $\pi'$. Observe that if $\hat\pi$ is a subproof of $\pi$ but not a proper subproof of $\pi_0$ then the non-terminals $\mathsf{N}_{\hat\pi}^j$ and $\mathsf{N}_{\hat\pi'}^j$ are of the same type for each $j$. Fix $i < |\Gamma|$ and a normal term $t_0 = t(\mathsf{N}_\pi^i/x)$. Suppose $t_0 \to t_1 \to \cdots \to t_k = r$ is a derivation in $\mathscr{H}$ of a $\Sigma$-term $r$. By Lemma 5.10, $t_i$ is normal for every $i \leq k$ and, without loss of generality, we may assume $t$ does not feature any non-terminals labelled by proofs with $\pi_0$ as a sub-proof. Throughout this derivation, recursively replace each occurrence of a non-terminal $\mathsf{N}_{\hat\pi}^j$ for which $\hat\pi$ is not a proper subproof of $\pi_0$ by the non-terminal $\mathsf{N}_{\hat\pi'}^j$. Arguing by induction on $k$, using $\pi_0 \sqsubset \pi_0'$, we deduce $t(\mathsf{N}_{\pi'}^i/x) \to^* s$ for some $\Sigma$-term $s$ with $s^\circ = r^\circ$. $\qquad\square$

We now turn our attention to the analysis of the subsumption relation with respect to the cut reduction and permutation steps of Figures 2 and 3. Only the most interesting cases will be covered in detail: the cut and quantifier permutation, and contraction and quantifier reduction. As before, we leave instances of the permutation inference implicit and make use of Lemma 5.19 without reference. Recall the characterisation of the $\mathsf{cut}$ inference from Remark 3.25:

$$\mathsf{cut}\ \frac{\pi_0 \vdash \Gamma, A \quad \pi_1 \vdash \Delta, \bar A}{\pi \vdash \Gamma, \Delta} \qquad \mathsf{N}_\pi^i a\vec x\vec y \sim \begin{cases} \mathsf{N}_{\pi_0}^i a\vec x(\mathsf{N}_{\pi_1}^n a\vec y(\mathsf{N}_{\pi_0}^m a\vec x)), & \text{if } u(A) > 0, \\ \mathsf{N}_{\pi_0}^i a\vec x(\mathsf{N}_{\pi_1}^n a\vec y), & \text{if } e(A) > 0, \\ \mathsf{N}_{\pi_0}^i a\vec x\langle\rangle, & \text{if } A \text{ is q.f.,} \end{cases}$$

where $i < |\Gamma|$ and $\vec x$, $\vec y$ and $a$ are terms of suitable type.

### 5.3.1 Cut Permutation

Suppose $\pi \rightsquigarrow \pi'$ are the two proofs



Due to the asymmetry in the production rules for cut, it is necessary to split the analysis of this reduction into two cases, depending on whether or not $A$ and $B$ are both universally quantified. Provided at least one of the two formulæ is existentially quantified or quantifier free, the two proofs above are equivalent and their languages are equal. This is proved in Lemma 5.22. If both $A$ and $B$ are universally quantified we do not expect equivalence to hold in general. However, if there are no contractions to the formula $\bar A$ in $\pi_1$ or the formula $\bar B$ in $\pi_2$, the proofs $\pi$ and $\pi'$ are equivalent. This is relevant to the cut reduction strategies employed in Theorem 1.1 and is treated in Lemma 5.25.

**Lemma 5.22.** *For $\pi \rightsquigarrow \pi'$ as above, if at least one of $u(A)$ and $u(B)$ is zero then $\pi \sim \pi'$.*

*Proof.* If one of $A$ or $B$ is quantifier-free the argument is straightforward following the production rules for cut. This leaves the following three cases to consider: $u(\bar A), u(B) >$

$0$, $u(A), u(\bar{B}) > 0$ and $u(\bar{A}), u(\bar{B}) > 0$. We consider only the first case as the second is symmetric and the third follows a simpler argument. Thus assume $e(A), u(B) > 0$.

Let $\vec{r}$, $\vec{s}$, $\vec{t}$ be sequences of terms of length $m = |\Gamma|$, $n = |\Delta|$ and $o = |\Lambda|$ respectively, and let $a : \varsigma$ be an arbitrary substitution stack. By the production rules for cut we have, for each $i \le m$, $j < n$ and $k < o$, and each term $w$, $w'$ of suitable type,

$$\mathsf{N}_{\hat{\pi}}^i a\vec{r}w\vec{s} \sim \begin{cases} \mathsf{N}_{\pi_0}^i a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s})w, & \text{if } i < m, \\ \mathsf{N}_{\pi_0}^{m+1} a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s})w, & \text{if } i = m, \end{cases} \qquad \mathsf{N}_{\hat{\pi}}^{m+1+j} a\vec{r}w\vec{s} \sim \mathsf{N}_{\pi_1}^j a\vec{s}(\mathsf{N}_{\pi_0}^m a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s})w)$$

$$\mathsf{N}_{\hat{\pi}'}^i a\vec{r}w'\vec{t} \sim \mathsf{N}_{\pi_0}^i a\vec{r}w'(\mathsf{N}_{\pi_2}^o a\vec{t}(\mathsf{N}_{\pi_0}^{m+1} a\vec{r}w')) \qquad \mathsf{N}_{\hat{\pi}'}^{m+1+k} a\vec{r}w'\vec{t} \sim \mathsf{N}_{\pi_2}^k a\vec{t}(\mathsf{N}_{\pi_0}^{m+1} a\vec{r}w')$$

In particular,

$$\hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{r}\vec{s} \sim \mathsf{N}_{\pi_0}^{m+1} a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s}) \tag{15}$$

and so, for $i \le m$,

$$\mathsf{N}_{\hat{\pi}'}^i a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s})\vec{t} \sim \mathsf{N}_{\pi_0}^i a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s})(\mathsf{N}_{\pi_2}^o a\vec{t}(\hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{r}\vec{s})). \tag{16}$$

We prove $\mathsf{N}_\pi^i a\vec{r}\vec{s}\vec{t} \sim \mathsf{N}_{\pi'}^i a\vec{r}\vec{s}\vec{t}$ for every $i < m + n + o$, from which Lemma 5.14 implies $\mathsf{N}_\pi^i \sim \mathsf{N}_{\pi'}^i$. For $i < m$ we have

$$\mathsf{N}_\pi^i a\vec{r}\vec{s}\vec{t} \sim \mathsf{N}_{\hat{\pi}}^i a\vec{r}(\mathsf{N}_{\pi_2}^o a\vec{t}(\hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{r}\vec{s}))\vec{s} \sim \mathsf{N}_{\pi_0}^i a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s})(\mathsf{N}_{\pi_2}^o a\vec{t}(\hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{r}\vec{s}))$$
$$\sim \mathsf{N}_{\hat{\pi}'}^i a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s})\vec{t}$$
$$\sim \mathsf{N}_{\pi'}^i a\vec{r}\vec{s}\vec{t}$$

by applying (16). For $j < n$,

$$\mathsf{N}_\pi^{m+j} a\vec{r}\vec{s}\vec{t} \sim \mathsf{N}_{\hat{\pi}}^{m+1+j} a\vec{r}(\mathsf{N}_{\pi_2}^o a\vec{t}(\hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{r}\vec{s}))\vec{s} \sim \mathsf{N}_{\pi_1}^j a\vec{s}(\mathsf{N}_{\pi_0}^m a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s})(\mathsf{N}_{\pi_2}^o a\vec{t}(\hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{r}\vec{s})))$$
$$\sim \mathsf{N}_{\pi_1}^j a\vec{s}(\mathsf{N}_{\hat{\pi}'}^m a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s})\vec{t})$$
$$\sim \mathsf{N}_{\pi'}^{m+j} a\vec{r}\vec{s}\vec{t}$$

again applying (16). For $k < o$, using (15):

$$\mathsf{N}_\pi^{m+n+k} a\vec{r}\vec{s}\vec{t} \sim \mathsf{N}_{\pi_2}^k a\vec{t}(\hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{r}\vec{s}) \sim \mathsf{N}_{\pi_2}^k a\vec{t}(\mathsf{N}_{\pi_0}^{m+1} a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s}))$$
$$\sim \mathsf{N}_{\hat{\pi}'}^{m+1+k} a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s})\vec{t}$$
$$\sim \mathsf{N}_{\pi'}^{m+n+k} a\vec{r}\vec{s}\vec{t}. \qquad \Box$$

As noted above, in the case $u(A)$ and $u(B)$ are both positive, language equality holds only in particular circumstances. A sufficient condition for this is given by the next lemma.

**Lemma 5.23.** *Let* $\pi \rightsquigarrow \pi'$ *be as above and assume* $u(A), u(B) > 0$. *Let* $\rho = \tau_{\hat{A}}^*$, $\sigma = \tau_{\hat{B}}^*$, $R = \hat{\mathsf{N}}_{\hat{\pi}'}^m a\vec{r}\vec{t}$ *and* $S = \hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{r}\vec{s}$. *If*

$$R \sim \hat{\mathsf{N}}_{\pi_0}^m a\vec{r}(\mathsf{N}_{\pi_2}^o a\vec{t}S) \qquad \text{and} \qquad S \sim \hat{\mathsf{N}}_{\pi_0}^{m+1} a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s}R)$$

*then* $\pi \sim \pi'$.

*Proof.* Recall that $R$ and $S$ have type $\rho$ and $\sigma$ respectively. Then for $i < m$,

$$\begin{aligned}
\mathsf{N}_\pi^i a\vec{r}\vec{s}\vec{t} &\sim \mathsf{N}_{\hat{\pi}}^i a\vec{r}(\mathsf{N}_{\pi_2}^o a\vec{t}S)\vec{s}\\
&\sim \mathsf{N}_{\pi_0}^i a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s}(\hat{\mathsf{N}}_{\pi_0}^m a\vec{r}(\mathsf{N}_{\pi_2}^o a\vec{t}S)))(\mathsf{N}_{\pi_2}^o a\vec{t}S)\\
&\sim \mathsf{N}_{\pi_0}^i a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s}R)(\mathsf{N}_{\pi_2}^o a\vec{t}(\hat{\mathsf{N}}_{\pi_0}^{m+1} a\vec{t}(\mathsf{N}_{\pi_1}^n a\vec{s}R)))\\
&\sim \mathsf{N}_{\hat{\pi}'}^i a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s}R)\vec{t}\\
&\sim \mathsf{N}_{\pi'}^i a\vec{r}\vec{s}\vec{t}
\end{aligned}$$

The other cases, namely $m \le i < n + o$ follow similar reasoning. $\qquad\square$

**Lemma 5.24.** *If $A, B \in \Pi_1 \cup \Sigma_1$ then $\pi \sim \pi'$.*

*Proof.* Assume $\bar{A}$ and $\bar{B}$ are both $\Sigma_1$ formulæ (if not, apply Lemma 5.22). Let $R$ and $S$ be as in Lemma 5.23. We have, by Lemma 5.17,

$$\begin{aligned}
Rw &\sim \mathsf{N}_{\pi_0}^m a\vec{r}w(\mathsf{N}_{\pi_2}^n a\vec{t}(\hat{\mathsf{N}}_{\pi_0}^{m+1} a\vec{r}w)) & Sw' &\sim \mathsf{N}_{\pi_0}^{m+1} a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s}(\hat{\mathsf{N}}_{\pi_0}^m a\vec{r}w'))w'\\
&\sim \mathsf{N}_{\pi_0}^m a\vec{r}w(\mathsf{N}_{\pi_2}^n a\vec{t}S) & &\sim \mathsf{N}_{\pi_0}^{m+1} a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s}R)w'
\end{aligned}$$

and hence

$$\begin{aligned}
R &\sim \hat{\mathsf{N}}_{\pi_0}^m a\vec{r}(\mathsf{N}_{\pi_2}^o a\vec{t}S) & S &\sim \hat{\mathsf{N}}_{\pi_0}^{m+1} a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s}R).
\end{aligned}$$

The previous lemma then implies $\pi \sim \pi'$. $\qquad\square$

**Lemma 5.25.** *For the same $\pi$ and $\pi'$, if there are no contractions to either the formula $\bar{A}$ in the sub-proof $\pi_1$ or the formula $\bar{B}$ in the sub-proof $\pi_2$ then $\pi \sim \pi'$.*

*Proof.* Suppose there are no contractions to $\bar{B}$ in $\pi_2$ and let $R$ and $S$ be as above. By Lemma 5.18, $\mathsf{N}_{\pi_2}^k a\vec{t}u \sim \mathsf{N}_{\pi_2}^k a\vec{t}v$ for any two terms $u, v : \tau_{\bar{B}}^*$. Hence, in particular,

$$\begin{aligned}
R &\sim \hat{\mathsf{N}}_{\pi_0}^m a\vec{r}(\mathsf{N}_{\pi_2}^o a\vec{t}S) & \mathsf{N}_{\pi_2}^o a\vec{t}S &\sim \hat{\mathsf{N}}_{\pi_0}^{m+1} a\vec{r}(\mathsf{N}_{\pi_1}^n a\vec{s}R)
\end{aligned}$$

which suffice, by the proof of Lemma 5.23, to show $\pi \sim \pi'$. $\qquad\square$

### 5.3.2 Contraction Reduction

Consider the two proofs



where $\pi_1^*$ denotes a copy of $\pi_1$ with fresh eigenvariables. Observe that $\pi_1 \sim \pi_1^*$.

Although the reduction above does not in general induce language inclusion, for the two scenarios required in Theorem 1.1, namely either $u(A) = 0$ or there are no applications of contraction are applied to the formula $\bar{A}$ in the sub-proof $\pi_1$, we have $\pi' \sqsubseteq \pi$. The following two lemmas deal with these two cases.

**Lemma 5.26.** *If $u(A) = 0$ then $\pi' \sqsubset \pi$.*

*Proof.* If $A$ is quantifier-free then $\pi \sim \pi'$ is easily established by following the reduction rules for cut. So assume $u(A) = 0 < u(\bar{A})$. Let $m = |\Gamma|$ and $n = |\Delta|$, and fix $i < m$ and $j < n$. Let $r = \mathsf{N}_{\pi_1}^n a\vec{y}$ and $r_* = \mathsf{N}_{\pi_1^*}^n a\vec{y}$. Unravelling the production rules for the two proofs yield

$$\mathsf{N}_\pi^i a\vec{x}\vec{y} \sim \mathsf{N}_{\pi_0}^i a\vec{x}rr \qquad \mathsf{N}_\pi^{m+j} a\vec{x}\vec{y} \sqsupset \{\mathsf{N}_{\pi_1}^j a\vec{y}(\mathsf{N}_{\pi_0}^{m+1} a\vec{x}rr), \mathsf{N}_{\pi_1}^j a\vec{y}(\mathsf{N}_{\pi_0}^m a\vec{x}rr)\}$$

$$\mathsf{N}_{\pi'}^i a\vec{x}\vec{y} \sim \mathsf{N}_{\pi_0}^i a\vec{x}r_*r \qquad \mathsf{N}_{\pi'}^{m+j} a\vec{x}\vec{y} \sim \{\mathsf{N}_{\pi_1}^j a\vec{y}(\mathsf{N}_{\pi_0}^{m+1} a\vec{x}r_*r), \mathsf{N}_{\pi_1^*}^j a\vec{y}(\mathsf{N}_{\pi_0}^m a\vec{x}r_*r)\}$$

Since $\pi_1 \sim \pi_1^*$, Lemma 5.13 implies $\pi' \sqsubset \pi$. $\qquad\qquad\square$

**Lemma 5.27.** *If $u(A) > 0$ and there are no contractions on the formula $\bar{A}$ in $\pi_1$, then $\pi' \sim \pi$.*

*Proof.* Suppose $u(A) > 0$. Let $\tau = \tau_{\bar{A}}^*$. Lemma 5.18 implies $\mathsf{N}_{\pi_1}^j a\vec{y}s \sim \mathsf{N}_{\pi_1}^j a\vec{y}\mathsf{c}_\tau$ for every $s : \tau$. Concerning derivations from $\pi$, this yields the following equivalences for $i < |\Gamma|$ and $j < |\Delta|$.

$$\mathsf{N}_\pi^i a\vec{x}\vec{y} \sim \mathsf{N}_{\hat{\pi}}^i a\vec{x}(\mathsf{N}_{\pi_1}^n a\vec{y}(\hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{x})) \qquad\qquad \mathsf{N}_\pi^{m+j} a\vec{x}\vec{y} \sim \mathsf{N}_{\pi_1}^j a\vec{y}(\hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{x})$$

$$\sim \mathsf{N}_{\hat{\pi}}^i a\vec{x}(\mathsf{N}_{\pi_1}^n a\vec{y}\mathsf{c}_\tau) \qquad\qquad\qquad\qquad \sim \mathsf{N}_{\pi_1}^j a\vec{y}\mathsf{c}_\tau$$

$$\sim \mathsf{N}_{\pi_0}^i a\vec{x}(\mathsf{N}_{\pi_1}^n a\vec{y}\mathsf{c}_\tau)(\mathsf{N}_{\pi_1}^n a\vec{y}\mathsf{c}_\tau)$$

Starting from $\pi'$ we obtain

$$\mathsf{N}_{\pi'}^i a\vec{x}\vec{y} \sim \mathsf{N}_{\hat{\pi}'}^i a\vec{x}(\mathsf{N}_{\pi_1^*}^n a\vec{y}(\hat{\mathsf{N}}_{\hat{\pi}}^m a\vec{x}\vec{y}))\vec{y} \qquad \mathsf{N}_{\pi'}^{m+j} a\vec{x}\vec{y} \sim \{\mathsf{N}_{\hat{\pi}'}^{m+1+j} a\vec{x}(\mathsf{N}_{\pi_1^*}^n a\vec{y}\mathsf{c}_\tau)\vec{y}, \mathsf{N}_{\pi_1^*}^j a\vec{y}\mathsf{c}_\tau\}$$

$$\sim \mathsf{N}_{\hat{\pi}'}^i a\vec{x}(\mathsf{N}_{\pi_1^*}^n a\vec{y}\mathsf{c}_\tau)\vec{y} \qquad\qquad\qquad\quad \sim \{\mathsf{N}_{\pi_1}^j a\vec{y}\mathsf{c}_\tau, \mathsf{N}_{\pi_1^*}^j a\vec{y}\mathsf{c}_\tau\}$$

$$\sim \mathsf{N}_{\pi_0}^i a\vec{x}(\mathsf{N}_{\pi_1^*}^n a\vec{y}\mathsf{c}_\tau)(\mathsf{N}_{\pi_1}^n a\vec{y}\mathsf{c}_\tau)$$

So $\pi' \sim \pi$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

### 5.3.3 Quantifier Permutation

Concerning permuting quantifier rules with cut, consider the following two proofs.

$$\forall_{\vec{\alpha}} \cfrac{\Gamma, A(\vec{\alpha}/\vec{v}), B}{\mathsf{cut} \cfrac{\hat{\pi} \vdash \Gamma, \forall \vec{v}A, B \qquad \Delta, \bar{B}}{\pi \vdash \Gamma, \forall \vec{v}A, \Delta}} \qquad \leadsto \qquad \forall \cfrac{\mathsf{cut} \cfrac{\Gamma, A(\vec{\alpha}/\vec{v}), B \qquad \Delta, \bar{B}}{\hat{\pi}' \vdash \Gamma, A(\vec{\alpha}/\vec{v}), \Delta}}{\pi' \vdash \Gamma, \forall \vec{v}A, \Delta} \qquad (17)$$

Let $\vec{\alpha} = (\alpha_j)_{j \leq p}$ and $\vec{v} = (v_j)_{j \leq p}$. Regularity ensures that $u(A) = 0$. In the following, if $\vec{u} = (u_j)_{j \leq p}$ is a sequence of terms of type $\iota$ and $u_{p+1} : \tau_A^*$, we write $\vec{u} \star u_{p+1}$ to abbreviate the sequence term $\langle u_0, \ldots, u_{p+1} \rangle : \tau_{\forall \vec{v}A}^*$.

Like with the case of permuting cuts, an application of the quantifier permutation reduction does not preserve equivalence of proofs in all cases. For the main theorem it suffice to prove only $\pi' \sqsubset \pi$. This is taken up in Lemma 5.29 below. First, however, we show that if $B$ is not universally quantified then indeed $\pi \sim \pi'$.

**Lemma 5.28.** *For $\pi$ and $\pi'$ above, if $u(B) = 0$ then $\pi' \sim \pi$.*

*Proof.* Suppose $u(B) = 0$ and $B$ is not quantifier-free. The other cases involve much similar arguments. Fix $\vec{r}$ and $\vec{s}$ sequences of normal terms of length $m = |\Gamma|$ and $n = |\Delta|$ respectively, and normal terms $t$ and $\vec{u} \star u'$ of type $\tau^*_{\forall \vec{v} A}$. By regularity of $\pi$ and Lemma 5.16,

$$\mathsf{N}^j_{\pi_1}([\vec{\alpha} \leftarrow \vec{u}]a)\vec{s} \sim \mathsf{N}^j_{\pi_1} a\vec{s}$$

for each $j \leq n$. Concerning $\pi$ the following equivalences therefore appear for $i \leq m$, $j < n$ and $k \leq m + 1$,

$$\mathsf{N}^i_\pi a\vec{r}t\vec{s} \sim \mathsf{N}^i_{\hat{\pi}} a\vec{r}t(\mathsf{N}^n_{\pi_1} a\vec{s}) \qquad\qquad \mathsf{N}^k_{\hat{\pi}} a\vec{r}(\vec{u} \star u') \sim \mathsf{N}^k_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}]a)\vec{r}u'$$
$$\mathsf{N}^{m+1+j}_\pi a\vec{r}t\vec{s} \sim \mathsf{N}^j_{\pi_1} a\vec{s}(\mathsf{N}^{m+1}_{\hat{\pi}} a\vec{r}t(\mathsf{N}^n_{\pi_1} a\vec{s}))$$

So, if $t$ is a normal term and $t \sim \{\vec{u}_0 \star u'_0, \ldots, \vec{u}_l \star u'_l\}$ is given by Lemma 5.7 then for $i \leq m$ and $j < n$,

$$\mathsf{N}^i_\pi a\vec{r}t\vec{s} \sim \{\mathsf{N}^i_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}_k]a)\vec{r}u'_k(\mathsf{N}^n_{\pi_1} a\vec{s}) \mid k \leq l\}$$
$$\mathsf{N}^{m+1+j}_\pi a\vec{r}t\vec{s} \sim \{\mathsf{N}^j_{\pi_1} a\vec{s}(\mathsf{N}^{m+1}_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}_k]a)\vec{r}u'_k(\mathsf{N}^n_{\pi_1} a\vec{s})) \mid k \leq l\}$$

Examining $\pi'$, we observe

$$\mathsf{N}^i_{\pi'} a\vec{r}t\vec{s} \sim \{\mathsf{N}^i_{\hat{\pi}'}([\vec{\alpha} \leftarrow \vec{u}_k]a)\vec{r}u'_k\vec{s} \mid k \leq l\}$$
$$\sim \{\mathsf{N}^i_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}_k]a)\vec{r}u'_k(\mathsf{N}^n_{\pi_1} a\vec{s}) \mid k \leq l\}$$
$$\mathsf{N}^{n+j}_{\pi'} a\vec{r}t\vec{s} \sim \{\mathsf{N}^j_{\pi_1} a\vec{s}(\mathsf{N}^{m+1}_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}_k]a)\vec{r}u'_k(\mathsf{N}^n_{\pi_1} a\vec{s})) \mid k \leq l\}$$

Hence $\mathsf{N}^i_{\pi'} \sim \mathsf{N}^i_\pi$ for every $i \leq m + n$ and so $\pi' \sim \pi$. $\qquad\square$

**Lemma 5.29.** *For $\pi$ and $\pi'$ as in (17), $\pi' \sqsubset \pi$.*

*Proof.* Fix $\vec{r}$ and $\vec{s}$ sequences of terms of length $m = |\Gamma|$ and $n = |\Delta|$ respectively. Let $\vec{u} \star u' : \tau^*_{\forall \vec{v} A}$. Suppose $u(B) > 0$ and $i \leq m$. The other cases have been considered earlier or involve similar but simpler arguments. As was observed earlier,

$$\mathsf{N}^j_{\pi_1}([\vec{\alpha} \leftarrow \vec{u}]a)\vec{s} \sim \mathsf{N}^j_{\pi_1} a\vec{s}$$

for each $j \leq n$. With respect to $\pi'$ the following equivalences therefore appear.

$$\mathsf{N}^i_{\pi'} a\vec{r}(\vec{u} \star u')\vec{s} \sim \mathsf{N}^i_{\hat{\pi}'}([\vec{\alpha} \leftarrow \vec{u}]a)\vec{r}u'\vec{s}$$
$$\sim \mathsf{N}^i_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}]a)\vec{r}u'(\mathsf{N}^n_{\pi_1} a\vec{s}(\mathsf{N}^{m+1}_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}]a)\vec{r}u'))$$

whereas the rules for $\pi$ yield, for arbitrary $t : \tau^*_{\forall \vec{v} A}$,

$$\mathsf{N}^i_\pi a\vec{r}t\vec{s} \sim \mathsf{N}^i_{\hat{\pi}} a\vec{r}t(\mathsf{N}^n_{\pi_1} a\vec{s}(\mathsf{N}^{m+1}_{\hat{\pi}} a\vec{r}t)) \qquad\qquad \mathsf{N}^i_{\hat{\pi}} a\vec{r}(\vec{u} \star u') \sim \mathsf{N}^i_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}]a)\vec{r}u'$$

If $t$ is a normal term and $t \sim \{\vec{u}_0 \star u'_0, \ldots, \vec{u}_l \star u'_l\}$ is given by Lemma 5.7 then for each $i \leq m$,

$$\mathsf{N}^i_{\hat{\pi}} a\vec{r}t \sim \{\mathsf{N}^i_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}_k]a)\vec{r}u'_k \mid k \leq l\}$$
$$\mathsf{N}^i_\pi a\vec{r}t\vec{s} \sim \{\mathsf{N}^i_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}_k]a)\vec{r}u'_k(\mathsf{N}^n_{\pi_1} a\vec{s}(\mathsf{N}^{m+1}_{\pi_0}([\vec{\alpha} \leftarrow \vec{u}_j]a)\vec{r}u'_j)) \mid k, j \leq l\} \qquad (18)$$

whereas, due to pattern-matching in the production rule for $\pi'$,

$$\mathsf{N}^i_{\pi'}a\vec{r}t\vec{s} \sim \{\mathsf{N}^i_{\pi_0}([\vec{\alpha} \hookleftarrow \vec{u}_k]a)\vec{r}u'_k(\mathsf{N}^n_{\pi_1}a\vec{s}(\mathsf{N}^{m+1}_{\pi_0}([\vec{\alpha} \hookleftarrow \vec{u}_k]a)\vec{r}u'_k)) \mid k \le l\} \tag{19}$$

Hence $\mathsf{N}^i_{\pi'} \sqsubseteq \mathsf{N}^i_{\pi}$ and $\pi' \sqsubseteq \pi$. $\qquad\square$

The contrast between equations (18) and (19) demonstrates why $\pi \sqsubseteq \pi'$ need not hold in general.

### 5.3.4 Quantifier Reduction

Consider the reduction

$$
\begin{array}{c}
\cfrac{
\overbrace{\pi_0} \quad\quad \overbrace{\pi_1}
}{}
\end{array}
\quad\leadsto\quad
\tag{20}
$$

$$
\forall_{\vec{\alpha}}\cfrac{\Gamma, A(\vec{\alpha}/\vec{v})}{\hat{\pi}_0 \vdash \Gamma, \forall\vec{v}A} \quad \exists_{\vec{s}}\cfrac{\Delta, \bar{A}(\vec{s}/\vec{v})}{\hat{\pi}_1 \vdash \Delta, \exists\vec{v}\bar{A}}
$$

$$\mathsf{cut}\cfrac{\phantom{XXXXXXXXXXXXXXXXX}}{\pi \vdash \Gamma, \Delta}$$

$$\mathsf{cut}\cfrac{\Gamma, A(\vec{s}/\vec{v}) \qquad \Delta, \bar{A}(\vec{s}/\vec{v})}{\pi' \vdash \Gamma, \Delta}$$

**Lemma 5.30.** *If $\pi \leadsto \pi'$ is the reduction above then $\pi \sim \pi'$.*

*Proof.* Let $m = |\Gamma|$, $n = |\Delta|$, $\vec{\alpha} = (\alpha_i)_{i \le p}$ and $\vec{s} = (s_i)_{i \le p}$. Recall that $\vec{s} \cdot a = (s_i \cdot a)_{i \le p}$. Note that regularity of $\pi$ implies $u(A) = 0$. This leaves two cases to consider: $A$ is quantifier-free or $e(A) > 0$. Suppose the latter, so the cut in $\pi'$ remains a quantified cut (the case $A$ is q.f. follows an analogous argument). The following equivalences arise, where $i < m$ and $j < n$.

$$
\begin{aligned}
\mathsf{N}^i_{\pi}a\vec{r}\vec{t} &\sim \mathsf{N}^i_{\hat{\pi}_0}a\vec{r}(\mathsf{N}^n_{\pi_1}a\vec{t}(\mathsf{N}^m_{\hat{\pi}_0}a\vec{r})) & \mathsf{N}^{m+j}_{\pi}a\vec{r}\vec{t} &\sim \mathsf{N}^j_{\hat{\pi}_1}a\vec{t}(\mathsf{N}^m_{\hat{\pi}_0}a\vec{r}) \\
&\sim \mathsf{N}^i_{\hat{\pi}_0}a\vec{r}(\vec{s} \cdot a \star \mathsf{N}^n_{\pi_1}a\vec{t}) & &\sim \mathsf{N}^j_{\pi_1}a\vec{t}(\mathsf{N}^m_{\hat{\pi}_0}a\vec{r}(\vec{s} \cdot a \star \mathsf{N}^n_{\pi_1}a\vec{t})) \\
&\sim \mathsf{N}^i_{\pi_0}([\vec{\alpha} \hookleftarrow \vec{s} \cdot a]a)\vec{r}(\mathsf{N}^n_{\pi_1}a\vec{t}) & &\sim \mathsf{N}^j_{\pi_1}a\vec{t}(\mathsf{N}^m_{\pi_0}([\vec{\alpha} \hookleftarrow \vec{s} \cdot a]a)\vec{r}(\mathsf{N}^n_{\pi_1}a\vec{t})) \\
&\sim \mathsf{N}^i_{\pi_0^{(\vec{s}/\vec{\alpha})}}a\vec{r}(\mathsf{N}^n_{\pi_1}a\vec{t}) & &\sim \mathsf{N}^j_{\pi_1}a\vec{t}(\mathsf{N}^m_{\pi_0^{(\vec{s}/\vec{\alpha})}}a\vec{r}(\mathsf{N}^n_{\pi_1}a\vec{t})) \\
&\sim \mathsf{N}^i_{\pi'}a\vec{r}\vec{t} & &\sim \mathsf{N}^{m+j}_{\pi'}a\vec{r}\vec{t}.
\end{aligned}
$$

The penultimate equivalence in each column is given by Lemma 5.16. $\qquad\square$

### 5.3.5 Remaining Reductions

The remaining rules are all straightforward to analyse and all induce language equality with the exception of weakening reduction for which only language inclusion holds in general.

## 5.4 Proof of Main Theorem

We can now prove Theorem 1.1. Let $\pi \vdash \exists\vec{v}F_{qf}$ be a regular proof and $\pi = \pi_0 \leadsto \pi_1 \leadsto \cdots \leadsto \pi_n$ be a reduction of $\pi$ to a quasi cut-free proof $\pi_n$ such that for each $i < n$, the reduction $\pi_i \leadsto \pi_{i+1}$ applies a cut reduction or permutation rule from Figures 2 or 3

to a sub-proof of $\pi_i$ with the restriction that a rule reducing the strong quantifier side of a cut is applied only if no other reduction of this cut is possible. By Lemma 5.21 and the analysis in the previous section, $L(\pi_{i+1}) \subseteq L(\pi_i)$ for each $i < n$. This together with Lemma 3.31 establishes part (iii) of the theorem. The existence of a reduction of the form above is well-known: see, e.g. [42], hence (i). Acyclicity of $\mathscr{H}_\pi$ is shown in Lemma 3.26, the bound on the order of $\mathscr{H}_\pi$ is given by Corollary 3.34, and the language bound in (ii) follows from Theorem 3.35.

## 6 Discussion

This work contributes to the structural analysis of first-order proofs with respect to their Herbrand content. To a first-order classical proof $\pi \vdash F$ of a $\Sigma_1$ formula we associate a recursion scheme $\mathscr{H}$ with a finite language that constitutes a Herbrand set for $F$. More generally, the language of $\mathscr{H}$ covers the Herbrand set implicit in any quasi cut-free proof obtained from $\pi$ by a sequence of reductions fulfilling the following two restrictions.

1. A contraction on a universally quantified formula is reduced only when no other reduction rule is applicable to this cut;

2. If two cuts are permuted in the following form then either there are no contractions on the formula $\bar{B}$ in the relevant subproof, or one of $A$ and $B$ is not universally quantified.

$$\text{cut } \frac{\dfrac{\Gamma, A, B \quad \Delta, \bar{A}}{\Gamma, \Delta, B} \quad \Lambda, \bar{B}}{\Gamma, \Delta, \Lambda} \quad \rightsquigarrow \quad \text{cut } \frac{\dfrac{\Gamma, A, B \quad \Lambda, \bar{B}}{\Gamma, \Lambda, A} \quad \Delta, \bar{A}}{\Gamma, \Delta, \Lambda}$$

The size of the Herbrand set is bounded by $2_{n+2}^{4|\pi|^3}$ where $|\pi|$ is the number of inferences in $\pi$ and $n$ is the maximal quantifier rank of a cut in $\pi$. Comparing with related work, the bound on the cardinality of the Herbrand expansion obtained by Gerhardy and Kohlenbach [20, Corollary 15] is $2_{\mathrm{dg}(\phi)+1}^{3\|t\|}$ where $\phi$ is a proof in Shoenfield's calculus [41], $t$ is the realiser extracted from $\phi$, and $\|t\|$ is the number of symbols in $t$. The degree $\mathrm{dg}(\phi)$ is the maximal $\neg$-depth of a cut formula in $\phi$. The $\neg$-depth of a formula is defined precisely in the discussion on pp. 17–25 of [19] as the maximal number of nested negations over quantifier-free sub-formulæ (that may contain an arbitrary number of negations). This is sufficient for describing the height of the tower of exponentials since, in Shoenfield's system, $\exists x$ is considered an abbreviation of $\neg \forall x \neg$. Thus (the translation of) a $\Pi_n \cup \Sigma_n$ formula has $\neg$-depth at most $n$. Presumably it is possible to give a polynomial translation from the sequent calculus into Shoenfield's system which preserves the maximal $\neg$-depth of cut formulæ (but, to the knowledge of the authors, this has not been done in the literature) and, moreover, to bound $\|t\|$ polynomially in terms of the number of inferences of $\phi$. Under these assumptions, the bound of Gerhardy and Kohlenbach would yield the upper bound $2_{n+1}^{p(|\pi|)}$ on the cardinality of a Herbrand expansion for some polynomial $p$ and any sequent calculus proof $\pi$ with $\Pi_n \cup \Sigma_n$-cuts. This would be one exponent less than our own.

Closely related is a bound obtained by Buss in [11]. The proof of Theorem 9 of [11] shows that, given a proof $\pi$ where all cut formulæ are contained in $\Pi_n \cup \Sigma_n$, there is a cut-free proof whose number of inferences is no greater than $2_{n+2}^{|\pi|}$. As an immediate corollary this also yields the upper bound of $2_{n+2}^{|\pi|}$ on the cardinality of the Herbrand expansion. If one is interested in the cardinality of the Herbrand expansion, Buss's bound and our Theorem 1.1 give the same number of iterations of the exponential function, but Gerhardy and Kohlenbach's would give one less. If one is interested in the number of inferences in the cut-free proof, Buss's bound is one exponential better than ours but has the same number of exponentials as the one that could be obtained from Gerhardy and Kohlenbach's since the number of inferences is at most exponential in the cardinality of the Herbrand expansion (considering the symbolic complexity of the end-sequent is constant). That being said, the bounds we obtain apply to any cut-free proof (and Herbrand disjunction) that can be reached by the class of reductions pertaining to 1 and 2 above. In particular, it places no restriction on which cut is to be reduced at any given step, and therefore accommodates a variety of strategies, including top-most and maximal cut-complexity. Whether this freedom of strategies necessitates the larger bound is not entirely clear, and requires further investigation.

Our approach provides a framework that appears well-suited for extensions. Below we highlight finer features of our representation of Herbrand's theorem and some potential applications.

**Sequent-based versus trace-based grammars**   In this paper, the grammar associated to a proof is 'sequent'-based in the following sense. Consider an inference of the form

$$\mathsf{r}\,\frac{G_0,\ldots,G_n}{F_0,\ldots,F_m}$$

The production rules corresponding to $\mathsf{r}$ can be seen as transforming a sequence of inputs $(x_0,\ldots,x_m)$ for the formulæ $F_0,\ldots,F_m$ to a sequence of terms $(t_0,\ldots,t_n)$ which are used as inputs for $G_0,\ldots,G_n$ in the inference rule immediately above $\mathsf{r}$. The production rules effect the whole sequence of inputs regardless of which formula is active. This is in contrast with the 'trace'-based grammars of, e.g., [24, 1, 2] where an inference of the form

$$\mathsf{r}\,\frac{\Gamma,G}{\Gamma,F}$$

is associated a production rule that updates an input for $F$ to an input for $G$, entirely ignoring presence of formulæ in $\Gamma$. In the latter type of grammars the derivations can be viewed as traces that climb up and also down the proof tree essentially mimicking the traces revealed through Gentzen-style cut-elimination. These grammars are generally cyclic and it is necessary to place equality constraints (the 'rigidity' conditions of [24, 1]) on derivations to ensure finite languages. For proofs that contain cuts with complexity greater than $\Pi_2/\Sigma_2$ the trace-based analysis quickly becomes infeasible. In contrast, the sequent-based approach generates an acyclic term grammar that not only ensures a finite language but allows one to obtain upper bounds on language size by standard language-theoretic arguments.

**Providing a minimal grammar** As mentioned in the introduction, part of the motivation behind this study is to ultimately invert the cut-elimination procedure and find an algorithmic method for introducing cuts into cut-free proofs. The idea has been successfully carried out for $\Pi_1/\Sigma_1$-cut introduction and more recently for the introduction of a single $\Pi_2/\Sigma_2$-cut (see [26, 25, 27, 32]). The general method proceeds as follows. Given a cut-free proof $\pi$, one first computes a concise representation of $\pi$ as a term grammar (such as a regular tree grammar whose language contains the Herbrand set induced by $\pi$). This grammar is then viewed as a proof with cut, in which the cut-formulæ are yet to be determined. Finding the cut-formulæ involves solving a unification problem induced by the grammar. Key to successfully carrying out this procedure is identifying natural classes of formal grammars that describe the instantiation structure of a proof with cut. Higher order recursion schemes provide a promising candidate to lift the method of cut-introduction above the $\Pi_2$ level.

**First-order logic in finite types** A natural extension to consider is first-order logic in finite (simple) types, namely many-sorted predicate logic with a sort of individuals for each simple type and well-typed application as a term forming operation. On the sequent calculus side, we add new quantifier inferences for each type:

$$\forall_\alpha^\sigma \frac{\Gamma, A(\alpha^\sigma/v^\sigma)}{\Gamma, \forall v^\sigma A} \qquad\qquad \exists_r^\sigma \frac{\Gamma, A(r^\sigma/v^\sigma)}{\Gamma, \exists v^\sigma A}$$

Here one can use the type hierarchy underpinning higher order recursion schemes. The formula types $\tau_A$ and $\tau_A^*$ are extended to incorporate higher-type quantification (for example $\tau_{\exists v^\sigma F} = \sigma \times \tau_F$ if $u(F) = 0$). The corresponding production rules for the new quantifier inferences will be identical to the rules for the ground type, though the move to higher-type means that substitution stacks can contain substitutions of higher-type symbols. We expect the analogous language preservation lemmas to hold. Moreover, Herbrand schemes for higher-type logic may provide a direct way to study the relation between the present work and the approach via functional interpretation in [20].

**Lifting the prenex restriction** Our representation of first-order proofs as recursion schemes forces an asymmetric interpretation of formulæ to types that does not easily generalise to non-prenex cuts. Specifically, the type of an existentially quantified formula is (except in the case of $\Sigma_1$) an order higher than the dual universally quantified formula. This disparity is due to the production rules for cut which treat the cut formula from one premise as a function which receives as its input 'witnesses' for the dual (cut) formula in the other premise. If the same representation is applied to non-prenex cuts then we must have that the type associated to a disjunction is an order higher than that assigned to the dual conjunction. Writing sound production rules for the disjunction and conjunction introduction inferences then becomes non-trivial.

One possible remedy is to switch to a two-sided sequent calculus. The order distinction forced by the production rules for cut becomes a distinction between the two sides of the sequent arrow thus permitting a more uniform association of formulæ to types that may permit a generalisation to cuts of any form. The hurdle here, however, will be in a satisfactory interpretation of negation which requires further investigation.

# 7 References

[1] Bahareh Afshari, Stefan Hetzl and Graham E. Leigh. 'Herbrand disjunctions, cut elimination and context-free tree grammars'. In: *13th International Conference on Typed Lambda Calculi and Applications (TLCA 2015)*. Ed. by Thorsten Altenkirch. Vol. 38. Leibniz International Proceedings in Informatics (LIPIcs). Dagstuhl, Germany: Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2015, pp. 1–16.

[2] Bahareh Afshari, Stefan Hetzl and Graham E. Leigh. 'Herbrand confluence for first-order proofs with $\Pi_2$-cuts'. In: *Proof: Concepts of Proof in Mathematics, Philosophy, and Computer Science*. Vol. 6. Ontos Mathematical Logic. De Gruyter, 2016, pp. 4–50.

[3] Bahareh Afshari, Stefan Hetzl and Graham E. Leigh. 'Herbrand's theorem revisited'. In: *Proceedings of Applied Mathematics and Mechanics* 16.1 (2016), pp. 905–906.

[4] Bahareh Afshari, Stefan Hetzl and Graham E. Leigh. 'On the Herbrand content of LK'. In: *Proceedings of Sixth International Workshop on Classical Logic and Computation (CL&C 2016)*. Ed. by Ulrich Kohlenbach, Steffen van Bakel and Stefano Berardi. Vol. 213. EPTCS. 2016, pp. 1–10.

[5] Jeremy Avigad and Solomon Feferman. 'Gödel's functional ("Dialectica") interpretation'. In: *The Handbook of Proof Theory*. Ed. by Sam Buss. North-Holland, 1999, pp. 337–405.

[6] Matthias Baaz and Stefan Hetzl. 'On the non-confluence of cut-elimination'. In: *Journal of Symbolic Logic* 76.1 (2011), pp. 313–340.

[7] Matthias Baaz, Stefan Hetzl et al. 'Cut-elimination: Experiments with CERES'. In: *Logic for Programming, Artificial Intelligence, and Reasoning (LPAR 2004)*. Ed. by Franz Baader and Andrei Voronkov. Vol. 3452. Lecture Notes in Computer Science. Springer, 2005, pp. 481–495.

[8] Matthias Baaz and Alexander Leitsch. 'Cut-elimination and redundancy-elimination by resolution'. In: *Journal of Symbolic Computation* 29.2 (2000), pp. 149–176.

[9] Franco Barbanera, Stefano Berardi and Massimo Schivalocchi. '"Classical" programming-with-proofs in $\lambda_{PA}^{Sym}$: An analysis of non-confluence'. In: *Theoretical Aspects of Computer Software*. Ed. by Martín Abadi and Takayasu Ito. Vol. 1281. Lecture Notes in Computer Science. Springer Berlin Heidelberg, 1997, pp. 365–390.

[10] Arnold Beckmann. 'Exact bounds for lengths of reductions in typed lambda-calculus.' In: *Journal of Symbolic Logic* 66.3 (2001), pp. 1277–1285.

[11] Sam Buss. 'Cut Elimination in Situ'. In: *Gentzen's Centenary: The Quest for Consistency*. Ed. by R. Kahle and M. Rathjen. Springer, 2015, pp. 245–277.

[12] Thierry Coquand. 'A semantics of evidence for classical arithmetic'. In: *Journal of Symbolic Logic* 60.1 (1995), pp. 325–337.

[13] Sebastian Eberhard and Stefan Hetzl. 'On the compressibility of finite languages and formal proofs'. To appear in *Information and Computation*.

[14]   Sebastian Eberhard and Stefan Hetzl. 'Compressibility of finite languages by grammars'. In: *Descriptional Complexity of Formal Systems (DCFS 2015)*. Ed. by Jeffrey Shallit and Alexander Okhotin. Vol. 9118. Lecture Notes in Computer Science. Springer, 2015, pp. 93–104.

[15]   Sebastian Eberhard and Stefan Hetzl. 'Inductive theorem proving based on tree grammars'. In: *Annals of Pure and Applied Logic* 166.6 (2015), pp. 665–700.

[16]   Gabriel Ebner, Stefan Hetzl et al. 'System Description: GAPT 2.0'. In: *8th International Joint Conference on Automated Reasoning (IJCAR 2016)*. Ed. by Nicola Olivetti and Ashish Tiwari. Vol. 9706. Lecture Notes in Computer Science. Springer, 2016, pp. 293–301.

[17]   Gabriel Ebner, Stefan Hetzl et al. 'On the generation of quantified lemmas'. Submitted. 2018.

[18]   Gerhard Gentzen. 'Untersuchungen über das logische Schließen II'. In: *Mathematische Zeitschrift* 39.3 (1935), pp. 405–431.

[19]   Philipp Gerhardy. 'Applications of Proof Interpretations'. PhD thesis. University of Aarhus, 2006.

[20]   Philipp Gerhardy and Ulrich Kohlenbach. 'Extracting Herbrand disjunctions by functional interpretation'. In: *Archive for Mathematical Logic* 44 (2005), pp. 633–644.

[21]   Kurt Gödel. 'Über eine noch nicht benützte Erweiterung des finiten Standpunktes'. In: *Dialectica* 12 (1958), pp. 280–287.

[22]   Willem Heijltjes. 'Classical proof forestry'. In: *Annals of Pure and Applied Logic* 161.11 (2010), pp. 1346–1366.

[23]   Hugo Herbelin. 'Séquents qu'on calcule'. PhD thesis. Université Paris 7, 1995.

[24]   Stefan Hetzl. 'Applying tree languages in proof theory'. In: *Language and Automata Theory and Applications (LATA 2012)*. Ed. by Adrian-Horia Dediu and Carlos Martín-Vide. Vol. 7183. Lecture Notes in Computer Science. Springer, 2012, pp. 301–312.

[25]   Stefan Hetzl, Alexander Leitsch, Giselle Reis and Daniel Weller. 'Algorithmic introduction of quantified cuts'. In: *Theoretical Computer Science* 549 (2014), pp. 1–16.

[26]   Stefan Hetzl, Alexander Leitsch and Daniel Weller. 'Towards algorithmic cut-introduction'. In: *Logic for Programming, Artificial Intelligence and Reasoning (LPAR 2012)*. Ed. by Nikolaj Bjørner and Andrei Voronkov. Vol. 7180. Lecture Notes in Computer Science. Springer, 2012, pp. 228–242.

[27]   Stefan Hetzl, Alexander Leitsch et al. 'Introducing quantified cuts in logic with equality'. In: *7th International Joint Conference on Automated Reasoning (IJCAR 2014)*. Ed. by Stéphane Demri, Deepak Kapur and Christoph Weidenbach. Vol. 8562. Lecture Notes in Computer Science. Springer, 2014, pp. 240–254.

[28] Stefan Hetzl and Lutz Straßburger. 'Herbrand-confluence for cut-elimination in classical first-order logic'. In: *Computer Science Logic (CSL 2012)*. Ed. by Patrick Cégielski and Arnaud Durand. Vol. 16. Leibniz International Proceedings in Informatics (LIPIcs). Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2012, pp. 320–334.

[29] Stefan Hetzl and Lutz Straßburger. 'Herbrand-Confluence'. In: *Logical Methods in Computer Science* 9.4 (2013).

[30] Stefan Hetzl and Daniel Weller. 'Expansion trees with cut'. Preprint available at `http://arxiv.org/abs/1308.0428`. 2013.

[31] David Hilbert and Paul Bernays. *Grundlagen der Mathematik II*. Springer, 1939.

[32] Alexander Leitsch and Michael Peter Lettmann. 'The problem of $\Pi_2$-cut-introduction'. In: *Theoretical Computer Science* 706 (2018), pp. 83–116.

[33] Richard McKinley. 'Proof nets for Herbrand's theorem'. In: *ACM Transactions on Computational Logic* 14.1 (2013), 5:1–5:31.

[34] Dale Miller. 'A compact representation of proofs'. In: *Studia Logica* 46.4 (1987), pp. 347–370.

[35] Georg Moser and Richard Zach. 'The epsilon calculus and Herbrand complexity'. In: *Studia Logica* 82.1 (2006), pp. 133–155.

[36] C.-H. Luke Ong and Steven J. Ramsay. 'Verifying higher-order functional programs with pattern-matching algebraic data types'. In: *Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL 2011)*. ACM, 2011, pp. 587–598.

[37] Luke Ong. 'Recursion schemes, collapsible pushdown automata and higher-order model checking'. In: *Language and Automata Theory and Applications (LATA 2013)*. Ed. by Adrian-Horia Dediu, Carlos Martín-Vide and Bianca Truthe. Springer, 2013, pp. 13–41.

[38] Diana Ratiu and Trifon Trifonov. 'Exploring the computational content of the infinite pigeonhole principle'. In: *Journal of Logic and Computation* 22.2 (2012), pp. 329–350.

[39] Helmut Schwichtenberg. 'Complexity of normalization in the pure typed lambda-calculus'. In: *The L.E.J.Brouwer Centenary Symposium*. Ed. by Anne S. Troelstra and D. van Dalen. North-Holland, 1982, pp. 453–457.

[40] Monika Seisenberger. 'On the constructive content of proofs'. PhD thesis. Ludwig-Maximilians-Universität München, 2003.

[41] Joseph R. Shoenfield. *Mathematical Logic*. 2nd edition. AK Peters, 2001.

[42] Anne S. Troelstra and Helmut Schwichtenberg. *Basic Proof Theory*. Cambridge: Cambridge University Press, 1996.

[43] Christian Urban. 'Classical logic and computation'. PhD thesis. University of Cambridge, 2000.