

---

# Contents

<b>Contents</b>	<b>1</b>
<b>1 Introduction</b>	<b>2</b>
1.1 Preliminaries . . . . .	2
1.2 Craig Interpolation . . . . .	3
<b>2 Calculi</b>	<b>4</b>
2.1 Resolution . . . . .	4
2.2 Resolution and Interpolation . . . . .	4
2.2.1 Interpolation and Skolemisation . . . . .	5
2.2.2 Interpolation and structure-preserving Normal Form Transformation	6
2.3 Sequent Calculus . . . . .	7
<b>3 Reduction to First-Order Logic without Equality</b>	<b>9</b>
3.1 Translation of formulas . . . . .	9
3.2 Computation of interpolants . . . . .	12
3.3 Proof by reduction . . . . .	16
<b>4 Proofs</b>	<b>18</b>
4.1 WT: Interpolation extraction in one pass . . . . .	18
4.2 WT: Interpolation extraction in two passes . . . . .	18
4.2.1 huang proof: propositional . . . . .	18
4.2.2 huang proof: overbinding . . . . .	22
4.2.3 final step of huang's proof . . . . .	27
<b>Bibliography</b>	<b>30</b>

# Introduction

## 1.1 Preliminaries

this section contains all the required notation but will just be written up nicely at the end

The language of a first-order formula  $A$  is denoted by  $L(A)$  and contains all predicate, constant and function symbols that occur in  $A$ . These are also referred to as the *non-logical symbols* of  $A$ . The *logical symbols* on the other hand include all logical connectives, quantifiers, the equality symbol ( $=$ ) as well as symbols denoting truth ( $\top$ ) and falsity ( $\perp$ ).

For formulas  $A_1, \dots, A_n$ ,  $L(A_1, \dots, A_n) = \bigcup_{1 \leq i \leq n} L(A_i)$ .

A term  $s$  is a subterm of a term  $t$  if  $s$  occurs in  $t$ .  $s$  is a strict subterm of  $t$  if  $s$  is a subterm of  $t$  and  $s \neq t$ .

An occurrence of  $\Phi$ -term is called *maximal* if it does not occur as subterm of another  $\Phi$ -term. An occurrence of a colored term  $t$  is a maximal colored term if it does not occur as subterm of another colored term.

We denote  $x_1, \dots, x_n$  by  $\bar{x}$ .

For a set of formulas  $\Phi$ ,  $\neg\Phi$  denotes  $\{\neg A \mid A \in \Phi\}$ .

A substitution is a mapping of variables to terms. It is denoted by  $\phi[x/t]$ , where  $\phi$  is a formula or term where each occurrence of the variable  $x$  is replaced by the term  $t$ . A substitution  $\sigma$  is called trivial on  $x$  if  $x\sigma = x$ . Otherwise it is called non-trivial.

An abstraction on the other hand is a mapping of terms to variables. It is denoted by  $\phi\{t/x\}$ , where  $\phi$  is a formula or term where each occurrence of the term  $t$  is replaced by the variable  $x$ .

The length of a term  $t$  is the number of symbols in  $t$ .

$A[s]_p$  denotes  $A$  with an occurrence of  $s$  at position  $p$ .

TODO: define  $A[s]$ .

## 1.2 Craig Interpolation

TODO: write some text about what interpolation means and that we prove more or less only reverse interpolation, but that's fine by the proposition

**Definition 1.1.** Let  $\Gamma$  and  $\Delta$  be sets of first-order formulas. An *interpolant* of  $\Gamma$  and  $\Delta$  is a first-order formula  $I$  such that

1.  $\Gamma \models I$
2.  $I \models \Delta$
3.  $L(I) \subseteq L(\Gamma) \cap L(\Delta)$ .

A *reverse interpolant* of  $\Gamma$  and  $\Delta$  is a first-order formula  $I$  such that  $I$  meets conditions 1 and 3 of an interpolant as well as:

$$2'. \Delta \models \neg I \quad \Delta$$

**Theorem 1.2** (Interpolation). *Let  $\Gamma$  and  $\Delta$  be sets of first-order formulas such that  $\Gamma \models \Delta$ . Then there exists an interpolant for  $\Gamma$  and  $\Delta$ .*

**Theorem 1.3** (Reverse Interpolation). *Let  $\Gamma$  and  $\Delta$  be sets of first-order formulas such that  $\Gamma \cup \Delta$  is unsatisfiable. Then there exists a reverse interpolant for  $\Gamma$  and  $\Delta$ .*

**Proposition 1.4.** *Theorem 1.2 and 1.3 are equivalent.*

*Proof.* Let  $\Gamma$  and  $\Delta$  be sets of first-order formulas such that  $\Gamma \models \Delta$ . Then  $\Gamma \cup \neg\Delta$  is unsatisfiable. By Theorem 1.3, there exists a reverse interpolant  $I$  for  $\Gamma$  and  $\neg\Delta$ . As  $\neg\Delta \models \neg I$ , we get by contraposition that  $I \models \Delta$ , hence  $I$  is an interpolant for  $\Gamma$  and  $\Delta$ .

For the other direction, let  $\Gamma$  and  $\Delta$  be sets of first-order formulas such that  $\Gamma \cup \Delta$  is unsatisfiable. Then  $\Gamma \models \neg\Delta$ , hence by Theorem 1.2, there exists an interpolant  $I$  of  $\Gamma$  and  $\neg\Delta$ . But as thus  $I \models \neg\Delta$ , we get by contraposition that  $\Delta \models \neg I$ , so  $I$  is a reverse interpolant for  $\Gamma$  and  $\Delta$ .  $\square$

As the notions of interpolation and reverse interpolation in this sense coincide, we will in the following only speak of interpolation where will be clear from the context which definition applies.

In the context of interpolation, every non-logical symbol is assigned a color which indicates its origin(s). A non-logical symbol is said to be  $\Gamma$  ( $\Delta$ )-*colored* if it only occurs in  $\Gamma$  ( $\Delta$ ) and *grey* in case it occurs in both  $\Gamma$  and  $\Delta$ . A symbol is *colored* if it is  $\Gamma$ - or  $\Delta$ -colored.

# Calculi

In this chapter, we introduce the calculi that are used subsequently. These are resolution and sequent calculus.

## 2.1 Resolution

Resolution calculus, in the formulation as given here, is a sound and complete calculus for first-order logic with equality. Due to the simplicity of its rules, it is widely used in the area of automated deduction.

**Definition 2.1.** A *clause* is a finite set of literals. The empty clause will be denoted by  $\square$ . A *resolution refutation* of a set of clauses  $\Gamma$  is a derivation of  $\square$  consisting of applications of resolution rules (cf. Figure 2.1) starting from clauses in  $\Gamma$ .  $\triangle$

**Theorem 2.2.** A clause set  $\Gamma$  is unsatisfiable if and only if there is resolution refutation of  $\Gamma$ .

*Proof.* See [Rob65].  $\square$

Clauses will usually be denoted by  $C$  or  $D$ , literals by  $l$  or  $l'$ , positions by  $p$ .

## 2.2 Resolution and Interpolation

In order to apply resolution to arbitrary first-order formulas, they have to be converted to clauses first. This usually makes use of intermediate normal forms which are defined as follows:

**Definition 2.3.** A formula is in *Negation Normal Form (NNF)* if negations only occur directly before of atoms. A formula is in *Conjunctive Normal Form (CNF)* if it is a conjunction of disjunctions of literals.  $\triangle$

In this context, the conjuncts of a CNF-formula are interpreted as clauses. A well-established procedure for the translation to CNF is comprised of the following steps:

$$\begin{aligned}
\text{Resolution: } & \frac{C \vee l \quad D \vee \neg l'}{(C \vee D)\sigma} \quad \sigma = \text{mgu}(l, l') \\
\text{Factorisation: } & \frac{C \vee l \vee l'}{(C \vee l)\sigma} \quad \sigma = \text{mgu}(l, l') \\
\text{Paramodulation: } & \frac{C \vee s = t \quad D[r]_p}{(C \vee D[t]_p)\sigma} \quad \sigma = \text{mgu}(s, r)
\end{aligned}$$

Figure 2.1: The rules of resolution calculus

1. NNF-Transformation
2. Skolemisation
3. CNF-Transformation

Step 1 can be achieved by solely pushing the negation inwards. As this transformation yields an equivalent formula, it clearly has no effect on the interpolants. Step 2 and 3 on the other hand do not produce equivalent formulas since they introduce new symbols. In this section, we will show that they nonetheless do preserve the set of interpolants. This fact is vital for the use of resolution-based methods for interpolant computation of arbitrary formulas.

### 2.2.1 Interpolation and Skolemisation

Skolemisation is a procedure for replacing existential quantifiers with Skolem terms:

**Definition 2.4.** Let  $V_{\exists x}$  be the set of universally bound variables in the scope of the occurrence of  $\exists x$  in a formula. The skolemisation of a formula  $A$  in NNF, denoted by  $\text{sk}(A)$ , is the result of replacing every occurrence of an existential quantifier  $\exists x$  in  $A$  by a term  $f(y_1, \dots, y_n)$  where  $f$  is a new Skolem function symbol and  $V_{\exists x} = \{y_1, \dots, y_n\}$ . In case  $V_{\exists x}$  is empty, the occurrence of  $\exists x$  is replaced by a new Skolem constant symbol  $c$ .

The skolemisation of a set of formulas  $\Phi$  is defined to be  $\text{sk}(\Phi) = \{\text{sk}(A) \mid A \in \Phi\}$ .  $\triangle$

**Proposition 2.5.** *Let  $\Gamma \cup \Delta$  be unsatisfiable. Then  $I$  is an interpolant for  $\Gamma \cup \Delta$  if and only if it is an interpolant for  $\text{sk}(\Gamma) \cup \text{sk}(\Delta)$ .*

*Proof.* Since  $\text{sk}(\cdot)$  adds fresh symbols to both  $\Gamma$  and  $\Delta$  individually, none of them are contained in  $L(\text{sk}(\Gamma)) \cap L(\text{sk}(\Delta))$ . Therefore condition 1.1.3 is satisfied in both directions.

As for any set of formulas  $\Phi$ , each model of  $\Phi$  can be extended to a model of  $\text{sk}(\Phi)$  and every model of  $\text{sk}(\Phi)$  is a witness for the satisfiability of  $\Phi$ ,  $\Phi \models I$  iff  $\text{sk}(\Phi) \models I$ . Hence conditions 1.1.1 and 1.1.2 remain satisfied for  $I$  as well.  $\square$

### 2.2.2 Interpolation and structure-preserving Normal Form Transformation

A common method for transforming a skolemised formula  $A$  into CNF while preserving their structure is defined as follows:

**Definition 2.6.** For every occurrence of a subformula  $B$  of  $A$ , introduce a new atom  $L_B$  which acts as a label for the subformula. For each of them, create a defining clause  $D_B$ :

If  $B$  is atomic:

$$D_B \equiv (\neg B \vee L_B) \wedge (B \vee \neg L_B)$$

If  $B$  is of the form  $\neg G$ :

$$D_B \equiv (L_B \vee L_G) \wedge (\neg L_B \vee \neg L_G)$$

If  $B$  is of the form  $G \wedge H$ :

$$D_B \equiv (\neg L_B \vee L_G) \wedge (\neg L_B \vee L_H) \wedge (L_B \vee \neg L_G \vee \neg L_H)$$

If  $B$  is of the form  $G \vee H$ :

$$D_B \equiv (L_B \vee \neg L_G) \wedge (L_B \vee \neg L_H) \wedge (\neg L_B \vee L_G \vee L_H)$$

If  $B$  is of the form  $G \supset H$ :

$$D_B \equiv (L_B \vee L_G) \wedge (L_B \vee \neg L_H) \wedge (\neg L_B \vee \neg L_G \vee L_H)$$

If  $B$  is of the form  $\forall xG$ :

$$D_B \equiv (\neg L_B \vee L_G) \wedge (L_B \vee \neg L_G)$$

Let  $\delta(A)$  be defined as  $\bigwedge_{B \in \Sigma(A)} D_B \wedge L_A$ , where  $\Sigma(A)$  denotes the set of occurrences of subformulas of  $A$ .  $\triangle$

Note that each of the  $D_B$  is in CNF, hence also  $\delta(A)$  for any skolemised formula  $A$ .

**Proposition 2.7.** *Let  $A$  be a formula. Then  $\text{sk}(A)$  is unsatisfiable if and only if  $\delta(\text{sk}(A))$  is unsatisfiable.*

**Proposition 2.8.** *Let  $\text{sk}(\Gamma) \cup \text{sk}(\Delta)$  be unsatisfiable. Then  $I$  is an interpolant for  $\text{sk}(\Gamma) \cup \text{sk}(\Delta)$  if and only if  $I$  is an interpolant for  $\delta(\text{sk}(\Gamma)) \cup \delta(\text{sk}(\Delta))$ .*

*Proof.* As  $\delta$  introduces fresh symbols for each  $\text{sk}(\Gamma)$  and  $\text{sk}(\Delta)$ , they must not occur in any interpolant for  $\text{sk}(\Gamma)$  and  $\text{sk}(\Delta)$ . This establishes condition 1.1.3 in both directions.

Using proposition 2.7, condition 1.1.1 and 1.1.2 are immediate.  $\square$

## 2.3 Sequent Calculus

The famous sequent calculus was introduced in [Gen35]. Its use of sequents in lieu of plain formulas allows for a natural mapping of the logical relations expressed by the connectives to the structure of proofs.

**Definition 2.9.** For multisets of first-order formulas  $\Gamma$  and  $\Delta$ ,  $\Gamma \vdash \Delta$  is called a *sequent*. In this context  $\Gamma$  forms the *antecedent*, whereas  $\Delta$  is referred to as *succedent*.

A sequent  $\Gamma \vdash \Delta$  is called *provable* if there is a sequent calculus proof of  $\Gamma \vdash \Delta$ .  $\triangle$

For the purposes of this thesis, we consider the cut-free fragment of sequent calculus.

**Theorem 2.10.** *Cut-free sequent calculus is sound and complete.*

The rules of cut-free sequent calculus are as follows:

### Axioms

$$A \vdash A$$

$$\vdash t = t$$

### Structural rules

- Contraction

$$\frac{\Gamma, A, A \vdash \Delta}{\Gamma, A \vdash \Delta} c : l$$

$$\frac{\Gamma \vdash \Delta, A, A}{\Gamma \vdash \Delta, A} c : r$$

- Weakening

$$\frac{\Gamma \vdash \Delta}{\Gamma, A \vdash \Delta} w : l$$

$$\frac{\Gamma \vdash \Delta}{\Gamma \vdash \Delta, A} w : r$$

### Propositional rules

- Negation

$$\frac{\Gamma \vdash \Delta, A}{\neg A, \Gamma \vdash \Delta} \neg : l$$

$$\frac{A, \Gamma \vdash \Delta}{\Gamma \vdash \Delta, \neg A} \neg : r$$

- Conjunction

$$\frac{\Gamma, A, B \vdash \Delta}{\Gamma, A \wedge B \vdash \Delta} \wedge : l$$

$$\frac{\Gamma \vdash \Delta, A \quad \Sigma \vdash \Pi, B}{\Gamma, \Sigma \vdash \Delta, \Pi, A \wedge B} \wedge : r$$

- Disjunction

$$\frac{\Gamma, A \vdash \Delta \quad \Sigma, B \vdash \Pi}{\Gamma, \Sigma, A \vee B \vdash \Delta, \Pi} \vee : l$$

$$\frac{\Gamma \vdash \Delta, A, B}{\Gamma \vdash \Delta, A \vee B} \vee : r$$

- Implication

$$\frac{\Gamma \vdash A, \Delta \quad \Sigma, B \vdash \Pi}{\Gamma, \Sigma A \supset B \vdash \Delta, \Pi} \supset : l$$

$$\frac{\Gamma, A \vdash \Delta, B}{\Gamma \vdash \Delta, A \supset B} \supset : r$$

### Quantifier rules

- Universal

$$\frac{\Gamma, A[x/t] \vdash \Delta}{\Gamma, \forall x A \vdash \Delta} \forall : l$$

$$\frac{\Gamma \vdash \Delta, A[x/y]}{\Gamma, \forall x A \vdash \Delta} \forall : r$$

- Existential

$$\frac{\Gamma, A[x/y] \vdash \Delta}{\Gamma, \exists x A \vdash \Delta} \exists : l$$

$$\frac{\Gamma \vdash \Delta, A[x/t]}{\Gamma, \exists x A \vdash \Delta} \exists : r$$

(provided no free variable of  $t$  becomes bound in  $A[x/t]$  and  $y$  does not occur free in  $\Gamma, \Delta$  or  $A$ )

### Equality rules

- Left rules

$$\frac{\Gamma, A[T/t] \vdash \Delta \quad \Sigma \vdash \Pi, s = t}{\Gamma, \Sigma, A[T/s] \vdash \Delta, \Pi} =: l_1 \quad \frac{\Gamma, A[T/s] \vdash \Delta \quad \Sigma \vdash \Pi, s = t}{\Gamma, \Sigma, A[T/t] \vdash \Delta, \Pi} =: l_2$$

- Right rules

$$\frac{\Gamma \vdash \Delta, A[T/t] \quad \Sigma \vdash \Pi, s = t}{\Gamma, \Sigma \vdash \Delta, \Pi, A[T/s]} =: r_1 \quad \frac{\Gamma \vdash \Delta, A[T/s] \quad \Sigma \vdash \Pi, s = t}{\Gamma, \Sigma \vdash \Delta, \Pi, A[T/t]} =: r_2$$

Figure 2.2: The rules of sequent calculus



# Reduction to First-Order Logic without Equality

A common theme of proofs in theoretical computer science is to avoid the tedious effort of proving the result from first principles by reducing the problem to one that is easier to solve. In this instance, we are able to give a reduction for finding interpolants in first-order logic *with* equality to first-order logic *without* equality, where it is simpler to give an appropriate algorithm.

The general layout of this approach is the following: From two sets  $\Gamma$  and  $\Delta$ , where  $\Gamma \cup \Delta$  is unsatisfiable, we compute two sets  $\Gamma'$  and  $\Delta'$  which do not make use of equality but simulate the effects of equality in  $\Gamma$  and  $\Delta$  via axioms. In the process of this transformation, also function symbols are replaced by predicate symbols with appropriate axioms to make sure that the behaviour of these function-representing predicates is compatible to the one of actual functions. Now an interpolant for  $\Gamma'$  and  $\Delta'$  can be derived using an algorithm that is only capable of handling predicate symbols as all other non-logical symbols have been removed. Since the additional axioms ensure that the newly added predicate symbols mimic equality and functions respectively, we will see that the occurrences of these predicates in the interpolant can be translated back to occurrences of equality and function symbols in first-order logic with equality in the language of  $\Gamma$  and  $\Delta$ , thereby yielding the originally desired interpolant.

## 3.1 Translation of formulas

As we shall see in this section, first-order formulas with equality can be transformed into first-order formulas without equality in a way that is satisfiability-preserving, which is sufficient for our purposes.

In order to simplify notation, we shall consider constant symbols to be function symbols of arity 0 in this section.

First, we define the axioms which allow for simulation of equality and functions in first order logic without equality and function symbols:

**Definition 3.1** (Equality and function axioms). For a first-order language  $\mathcal{L}$  and fresh predicate symbols  $E$  and  $F_f$  for  $f \in \text{FS}(\mathcal{L})$ , we define:

$$\begin{aligned} F_{\text{Ax}}(\mathcal{L}) &\stackrel{\text{def}}{=} \bigwedge_{f \in \text{FS}(\mathcal{L})} \forall \bar{x} \exists y (F_f(\bar{x}, y) \wedge (\forall z (F_f(\bar{x}, z) \supset E(y, z)))) \\ \text{Refl}(P) &\stackrel{\text{def}}{=} \forall x P(x, x) \\ \text{Congr}(P) &\stackrel{\text{def}}{=} \forall x_1 \forall y_1 \dots \forall x_{\text{ar}(P)} \forall y_{\text{ar}(P)} ((E(x_1, y_1) \wedge \dots \wedge E(x_{\text{ar}(P)}, y_{\text{ar}(P)})) \supset \\ &\quad (P(x_1, \dots, x_{\text{ar}(P)}) \supset P(y_1, \dots, y_{\text{ar}(P)}))) \\ E_{\text{Ax}}(\mathcal{L}) &\stackrel{\text{def}}{=} \text{Refl}(E) \wedge \bigwedge_{\substack{P \in \text{PS}(\mathcal{L}) \cup \{E\} \cup \\ \{F_f \mid f \in \text{FS}(\mathcal{L})\}}} \text{Congr}(P) \end{aligned} \quad \triangle$$

$\text{Refl}(P)$  will be referred to as reflexivity axiom of  $P$ ,  $\text{Congr}(P)$  as congruence axiom of  $P$ . Now we define the precise language we translate formulas to as well as the translation procedure.

**Definition 3.2** (Translation of languages). Let  $\mathcal{L}$  be a language. Then  $\text{T}(\mathcal{L})$  denotes  $(\text{L}(\mathcal{L}) \cup \{E\} \cup \{F_f \mid f \in \text{FS}(\mathcal{L})\}) \setminus (\{=\} \cup \text{FS}(\mathcal{L}))$ .  $\triangle$

**Definition 3.3** (Translation and inverse translation of formulas). Let  $A$  be a first-order formula and  $E$  and  $F_f$  for  $f \in \text{FS}(A)$  be fresh predicate symbols. Then  $\text{T}(A)$  is the result of applying the following algorithm to  $A$ :

1. Replace every occurrence of  $s = t$  in  $A$  by  $E(s, t)$
2. As long as there is an occurrence of a function symbol  $f$  in  $A$ :  
 Let  $B$  be the atom in which  $f$  occurs as outermost symbol of a term. Then  $B$  is of the form  $P(s_1, \dots, s_{j-1}, f(\bar{t}), s_{j+1}, \dots, s_m)$ . Replace  $B$  in  $A$  by  $\exists y (F_f(\bar{t}, y) \wedge P(s_1, \dots, s_{j-1}, y, s_{j+1}, \dots, s_m))$  for a fresh variable  $y$ .

Moreover, let the inverse operation  $\text{T}^{-1}(B)$  for formulas  $B$  in the language  $\text{T}(L(A))$  be defined as the result of applying the following algorithm to  $B$ :

1. Replace every occurrence of  $E(s, t)$  in  $B$  by  $s = t$ .
2. For every  $f \in \text{FS}(A)$ , replace every occurrence of  $\exists y (F_f(\bar{t}, y) \wedge P(s_1, \dots, s_{j-1}, y, s_{j+1}, \dots, s_m))$  in  $B$  by  $P(s_1, \dots, s_{j-1}, f(\bar{t}), s_{j+1}, \dots, s_m)$  and every remaining occurrence of  $F_f(\bar{t}, s)$  by  $f(\bar{t}) = s$ .

For sets of first-order formulas  $\Phi$ , let  $\text{T}(\Phi) \stackrel{\text{def}}{=} \bigcup_{A \in \Phi} \text{T}(A)$  and  $\text{T}^{-1}(\Phi) \stackrel{\text{def}}{=} \bigcup_{A \in \Phi} \text{T}^{-1}(A)$ .  $\triangle$

**Lemma 3.4.** Let  $A$  be a first-order formula and  $\Phi$  be a set of first-order formulas. Then  $\text{T}^{-1}(\text{T}(A)) = A$  and  $\text{T}^{-1}(\text{T}(\Phi)) = \Phi$ .

*Proof.* Step 1 and 2 in the transformation algorithm for  $\text{T}$  and  $\text{T}^{-1}$  are each concerned with a different set of symbols and therefore do not interfere with each other. Moreover, the respective steps in both algorithms are the inverse of each other. For step 1, this is immediate and for step 2, consider that all occurrences of  $F_f$  for  $f \in \text{FS}(A)$

in  $T(A)$  have been introduced by  $T$  and are consequently of the form  $\exists y(F_f(\bar{t}, y) \wedge P(s_1, \dots, s_{j-1}, y, s_{j+1}, \dots, s_m))$ , which is replaced by  $P(s_1, \dots, s_{j-1}, f(\bar{t}), s_{j+1}, \dots, s_m)$  by  $T^{-1}$ .  $\square$

**Definition 3.5** (Translation of formulas including axioms). For first-order formulas  $A$ , let  $T_{Ax}(A) = F_{Ax}(L(A)) \wedge E_{Ax}(L(A)) \wedge T(A)$  and for sets of first-order formulas  $\Phi$ , let  $T_{Ax}(\Phi) = \{F_{Ax}(L(\Phi)), E_{Ax}(L(\Phi))\} \cup T(\Phi)$ .  $\triangle$

Note that  $T_{Ax}(A)$  contains neither the equality predicate nor function symbols but additional predicate symbols instead. More formally:

**Lemma 3.6.**

1. Let  $\Phi$  be a set of first-order formulas. Then  $T_{Ax}(\Phi)$  is in the language  $T(L(\Phi))$ .
2. If  $\Psi$  is in the language  $T(\mathcal{L})$ , then  $T^{-1}(\Psi)$  is in the language  $\mathcal{L}$ .

**Proposition 3.7.** Let  $\Phi$  be a set of first-order formulas.

1. If  $\Phi$  is satisfiable, then so is  $T_{Ax}(\Phi)$ .
2. Let  $\mathcal{L}$  be a first-order language and  $\Phi$  a set of first-order formulas in the language  $T(\mathcal{L})$ . If  $\{F_{Ax}(\mathcal{L}), E_{Ax}(\mathcal{L})\} \cup \Phi$  is satisfiable, then so is  $T^{-1}(\Phi)$ .

*Proof.* Suppose  $\Phi$  is satisfiable. Let  $M$  be a model of  $\Phi$ . We show that  $T_{Ax}(\Phi)$  is satisfiable by extending  $M$  to the language  $L(\Phi) \cup \{E\} \cup \{F_f \mid f \in FS(A)\}$  and proving that the extended model satisfies  $T_{Ax}(\Phi)$ .

First, let  $M \models E(s, t)$  if and only if  $M \models s = t$ . By reflexivity of equality, it follows that  $M \models \text{Refl}(E)$ . As any predicate, in particular  $E$  and  $F_f$  for every  $f \in FS(\Phi)$ , satisfy the congruence axiom with respect to  $=$ , by the definition of  $E$  in  $M$ , they satisfy the congruence axiom with respect to  $E$ . Therefore  $M$  is a model of  $E_{Ax}(L(\Phi))$ .

Second, let  $M \models F_f(\bar{x}, y)$  if and only if  $M \models f(\bar{x}) = y$  for all  $f \in FS(\Phi)$ . Since  $M$  is a model of  $\Phi$ , it maps every function symbol  $f$  to a function, which by definition returns a unique result for every combination of parameters. This however is precisely the logical requirement on  $F_f$  stated by  $F_{Ax}(L(\Phi))$ , hence  $M$  is a model of  $F_{Ax}(L(\Phi))$ .

Lastly, we show that  $M \models T(A)$  for all  $A \in \Phi$ . By the above definition of  $E$  in  $M$ , step 1 of the algorithm in definition 3.3 yields a formula that is satisfied by  $M$  as it satisfies every formula of  $\Phi$ . For step 2, suppose  $P(s_1, \dots, s_{j-1}, f(\bar{t}), s_{j+1}, \dots, s_m)$  does (not) hold under  $M$ . Let  $y$  such that  $M \models f(\bar{t}) = y$ . By our definition of  $F$  under  $M$ ,  $M \models F(\bar{t}, y)$  with this unique  $y$ . Hence  $\exists y(F(\bar{t}, y) \wedge P(s_1, \dots, s_{j-1}, y, s_{j+1}, \dots, s_m))$  does (not) hold under  $M$ .

For the other direction, suppose  $\{F_{Ax}(\mathcal{L}), E_{Ax}(\mathcal{L})\} \cup \Phi$  is satisfiable. We extend a model  $M$  of this set of formulas to a model of  $T^{-1}(\Phi)$  by extending it from the language  $T(\mathcal{L})$  to include  $\{=\}$  and  $FS(\mathcal{L})$ .

First, let  $M \models s = t$  if and only if  $M \models E(s, t)$ . As  $M$  is a model of  $E_{Ax}(\mathcal{L})$ ,  $E$  is reflexive. Since  $M \models \text{Congr}(E)$ ,  $M \models \forall x \forall y (E(x, y) \wedge E(x, x) \supset (E(x, x) \supset E(y, x)))$ . As we know that  $E$  is reflexive, this simplifies to  $M \models \forall x \forall y (E(x, y) \supset E(y, x))$ , i.e.  $E$  is symmetric in  $M$ . We show the transitivity of  $E$  by another instance of  $\text{Congr}(E)$ :  $M \models \forall x \forall y \forall z ((E(y, x) \wedge E(y, z)) \supset (E(y, y) \supset E(x, z)))$ . As  $E$  is reflexive and symmetric,

we get that  $M \models \forall x \forall y \forall z ((E(x, y) \wedge E(y, z)) \supset E(x, z))$ . As these properties directly also apply to  $=$  in  $M$ , equality adheres to the required axioms in  $M$ .

Second, let  $M \models f(\bar{t}) = s$  if and only if  $M \models F_f(\bar{t}, s)$  for all  $f \in \text{FS}(\mathcal{L})$ . As by assumption  $M$  is a model of  $F_{\text{Ax}}(A)$ , we know that for every  $\bar{t}$ , some  $s$  with  $M \models F(\bar{t}, s)$  exists and is uniquely defined. Hence  $f$  in  $M$  refers to a well-defined function.

Lastly, to show that  $M \models T^{-1}(\Phi)$ , consider that the interpretations of the predicates  $E$  and  $=$  coincide in  $M$ . Furthermore, let  $B$  be an occurrence of  $\exists y (F_f(\bar{t}, y) \wedge P(s_1, \dots, s_{j-1}, y, s_{j+1}, \dots, s_m))$  for some  $f \in \text{FS}(\mathcal{L})$  in  $\Phi$ . Then by the above definition of  $f$  in  $M$ , we have that  $B$  is in  $M$  equivalent to  $\exists y f(\bar{t}) = y \wedge P(s_1, \dots, s_{j-1}, y, s_{j+1}, \dots, s_m)$ , which due to  $f$  being a function is equivalent to  $M \models P(s_1, \dots, s_{j-1}, f(\bar{t}), s_{j+1}, \dots, s_m)$ .

Similarly, let  $B$  be an occurrence of  $F_f(\bar{t}, s)$  in  $\Phi$ . Then by our above definition of  $f$  in  $M$ , we have that  $M \models f(\bar{t}) = s$  iff  $M \models B$ .  $\square$

**Corollary 3.8.** *Let  $\Phi$  be a set of first-order formulas. Then  $\Phi$  is satisfiable if and only if  $T_{\text{Ax}}(\Phi)$  is satisfiable.*

*Proof.* The left-to-right direction is directly given in Proposition 3.7. For the other direction, consider that by Proposition 3.7,  $T^{-1}(T(\Phi))$  is satisfiable, which by Lemma 3.4 is nothing else than  $\Phi$ .  $\square$

## 3.2 Computation of interpolants

For the proof of the interpolation theorem by reduction we require an algorithm that operates in first-order logic without equality and function symbols, which we describe in this section.

*Remark.* As the idea of this reduction is to simplify the problem by amongst others not considering function symbols, resolution-based methods can not be employed in a direct manner. This is because function symbols appear naturally in them as they usually handle existential quantification by means of skolemisation, i.e. a new function symbol is introduced for every occurrence of an existential quantifier in the scope of a universal quantifier. Translating the skolemised formulas to a language without function symbols as described in Definition 3.3 is of no avail since this translation introduces new existential quantifiers for every function symbol it encounters, necessitating skolemisation yet again.

We now show that interpolants can be computed by means of a sequent calculus based procedure by Maehara. It is slightly stronger than the required statement as it allows for interpolants of partitions of sequents:

**Definition 3.9** (Partition of sequents). A partition of a sequent  $\Gamma \vdash \Delta$  is denoted by  $\langle (\Gamma_1; \Delta_1), (\Gamma_2; \Delta_2) \rangle$ , where  $\Phi_1 \cap \Phi_2 = \emptyset$  and  $\Phi_1 \cup \Phi_2 = \Phi$  for  $\Phi \in \{\Gamma, \Delta\}$ .  $\triangle$

Note that while for partitions of sets it is usually required that the parts are non-empty, empty parts in a partition of sequents are permitted.

**Lemma 3.10** (Maehara). *Let  $\Gamma$  and  $\Delta$  be sets of first-order clauses without equality and function symbols such that  $\Gamma \vdash \Delta$  is provable in sequent calculus. Then for any partition  $\langle(\Gamma_1; \Delta_1), (\Gamma_2; \Delta_2)\rangle$  there is an interpolant  $I$  such that*

1.  $\Gamma_1 \vdash \Delta_1, I$  is provable
2.  $\Gamma_2, I \vdash \Delta_2$  is provable
3.  $L(I) \subseteq L(\Gamma_1, \Delta_1) \cap L(\Gamma_2, \Delta_2)$

*Proof.* We prove this lemma by induction on the number of inferences in a proof of  $\Gamma \vdash \Delta$ . As many cases are similar, we prove some examples only.

Base case. Suppose no rules were applied. Then  $C \vdash D$  is of one of the following forms:

- $C \vdash D$  is of the form  $A \vdash A$ . We give interpolants for any of the four possible partitions:
  1.  $\langle(A; A), (;)\rangle: I = \perp$
  2.  $\langle(;), (A; A)\rangle: I = \top$
  3.  $\langle(; A), (A;)\rangle: I = \neg A$
  4.  $\langle(A;), (; A)\rangle: I = A$
- $C \vdash D$  is of the form  $\vdash t = t$ . We again give interpolants for all possible partitions:
  1.  $\langle(; t = t), (;)\rangle: I = \perp$
  2.  $\langle(;), (; t = t)\rangle: I = \top$

Structural rules. Suppose the property holds for  $n$  rule applications and the  $(n + 1)$ th rule is a structural one. We consider some examples:

- The last rule application is an instance of  $c : l$ . Then it is of the form:

$$\frac{\Gamma, A, A \vdash \Delta}{\Gamma, A \vdash \Delta} c : l$$

There are two possible partition schemes: of  $\Gamma, A \vdash \Delta$ :

1.  $\chi = \langle(\Gamma_1, A; \Delta_1), (\Gamma_2; \Delta_2)\rangle$ . By the induction hypothesis, we know that there is an interpolant  $I$  for the partition  $\langle(\Gamma_1, A, A; \Delta_1), (\Gamma_2; \Delta_2)\rangle$  of the upper sequent.  $I$  serves as interpolant for  $\chi$  as well.
  2.  $\chi = \langle(\Gamma_1; \Delta_1), (\Gamma_2, A; \Delta_2)\rangle$ . By a similar argument, we get that there is an interpolant  $I$  for  $\langle(\Gamma_1; \Delta_1), (\Gamma_2, A, A; \Delta_2)\rangle$ , which again is also an interpolant for  $\chi$ .
- The last rule application is an instance of  $w : r$ . Then it is of the form:

$$\frac{\Gamma \vdash \Delta}{\Gamma \vdash \Delta, A} w : r$$

By the induction hypothesis, there exists an interpolant  $I$  for any partition  $\langle(\Gamma_1; \Delta_1), (\Gamma_2; \Delta_2)\rangle$  of  $\Gamma \vdash \Delta$ . Clearly  $I$  remains an interpolant when adding  $A$  to either  $\Delta_1$  or  $\Delta_2$ .

Propositional rules. Suppose the property holds for  $n$  rule applications and the  $(n+1)$ th rule is a propositional one. We consider some examples:

- The last rule application is an instance of  $\neg : l$ . Then it is of the form:

$$\frac{\Gamma \vdash \Delta, A}{\neg A, \Gamma \vdash \Delta} \neg : l$$

There are two possible partition schemes of  $\Gamma, \neg A \vdash \Delta$ :

1.  $\chi = \langle (\Gamma_1, \neg A; \Delta_1), (\Gamma_2; \Delta_2) \rangle$ . By the induction hypothesis, there exists an interpolant  $I$  for the partition  $\langle (\Gamma_1; \Delta_1, A), (\Gamma_2; \Delta_2) \rangle$  of the upper sequent. Clearly  $I$  is an interpolant for  $\chi$  as well.
  2.  $\chi = \langle (\Gamma_1; \Delta_1), (\Gamma_2, \neg A; \Delta_2) \rangle$ . A similar argument goes through.
- The last rule application is an instance of  $\supset : l$ . Then it is of the form:

$$\frac{\Gamma \vdash \Delta, A \quad \Sigma, B \vdash \Pi}{\Gamma, \Sigma, A \supset B \vdash \Delta, \Pi} \supset : l$$

There are two possible partition schemes of  $\Gamma, A \supset B \vdash \Delta$ :

1.  $\chi = \langle (\Gamma_1, \Sigma_1, A \supset B; \Delta_1, \Pi_1), (\Gamma_2, \Sigma_2; \Delta_2, \Pi_2) \rangle$ . By the induction hypothesis, there is an interpolant  $I_1$  for the partition  $\langle (\Gamma_1; \Delta_1, A), (\Gamma_2; \Delta_2) \rangle$  of the left upper sequent. Hence for  $I_1$ , we have that  $\Gamma_1 \vdash \Delta_1, A, I_1$  and  $I_1, \Gamma_2 \vdash \Delta_2$  are provable.

Moreover, we also get by the induction hypothesis that there is an interpolant  $I_2$  for the partition  $\langle (\Sigma_1, B; \Pi_1), (\Sigma_2; \Pi_2) \rangle$  of the right upper sequent. Therefore  $\Sigma_1, B \vdash \Pi_1, I_2$  and  $I_2, \Sigma_2 \vdash \Pi_2$  are provable.

Using these prerequisites, we first establish that  $I_1 \vee I_2$  fulfills conditions 1 and 2 of an interpolant for  $\chi$ :

$$\frac{\frac{\Gamma_1 \vdash \Delta_1, A, I_1 \quad \Sigma_1, B \vdash \Pi_1, I_2}{\Gamma_1, \Sigma_1, A \supset B \vdash \Delta_1, \Pi_1, I_1, I_2} \supset : l}{\Gamma_1, \Sigma_1, A \supset B \vdash \Delta_1, \Pi_1, I_1 \vee I_2} \vee : r$$

$$\frac{I_1, \Gamma_2 \vdash \Delta_2 \quad I_2, \Sigma_2 \vdash \Pi_2}{I_1 \vee I_2, \Gamma_2, \Sigma_2 \vdash \Delta_2, \Pi_2} \vee : l$$

To show that also condition 3 is satisfied, consider that by the induction hypothesis, it holds that:

$$L(I_1) \subseteq L(\Gamma_1, \Delta_1, A) \cap L(\Gamma_2, \Delta_2)$$

$$L(I_2) \subseteq L(\Sigma_1, B, \Pi_1) \cap L(\Sigma_2, \Pi_2)$$

Therefore

$$\begin{aligned}
 L(I_1) \cup L(I_2) &\subseteq (L(\Gamma_1, \Delta_1, A) \cap L(\Gamma_2, \Delta_2)) \cup (L(\Sigma_1, B, \Pi_1) \cap L(\Sigma_2, \Pi_2)) \\
 &\Downarrow \\
 L(I_1) \cup L(I_2) &\subseteq (L(\Gamma_1, \Delta_1, A) \cup L(\Sigma_1, B, \Pi_1)) \cap (L(\Gamma_2, \Delta_2) \cup L(\Sigma_2, \Pi_2)) \\
 &\Updownarrow \\
 L(I_1 \vee I_2) &\subseteq L(\Gamma_1, \Sigma_1, A \supset B, \Delta_1, \Pi_1) \cap L(\Gamma_2, \Sigma_2, \Delta_2, \Pi_2)
 \end{aligned}$$

2.  $\chi = \langle (\Gamma_1, \Sigma_1; \Delta_1, \Pi_1), (\Gamma_2, \Sigma_2, A \supset B; \Delta_2, \Pi_2) \rangle$ . The argument for this case is similar using  $I_1 \wedge I_2$  as interpolant.

Quantifier rules. Suppose the property holds for  $n$  rule applications and the  $(n + 1)$ th rule is a quantifier rule. We consider some examples:

- The last rule application is an instance of  $\forall : l$ . Then it is of the form:

$$\frac{\Gamma, A[x/t] \vdash \Delta}{\Gamma, \forall x A \vdash \Delta} \forall : l$$

where no free variable of  $t$  becomes bound in  $A[x/t]$ .

There are two possible partition schemes of  $\Gamma, \forall x A \vdash \Delta$ :

1.  $\langle (\Gamma_1, \forall x A; \Delta_1), (\Gamma_2; \Delta_2) \rangle$ . By the induction hypothesis, there is an interpolant  $I$  of the partition  $\langle (\Gamma_1, A[x/t]; \Delta_1), (\Gamma_2; \Delta_2) \rangle$ . Hence for  $I$ ,  $\Gamma_1, A[x/t] \vdash \Delta_1, I$  and  $I, \Gamma_2 \vdash \Delta_2$  are provable. By an application of  $\forall : l$  to the first sequent we get  $\Gamma_1, \forall x A \vdash \Delta_1, I$ , so  $I$  satisfies conditions 1 and 2 of being an interpolant for  $\chi$ .

In order to show that also  $L(I) \subseteq L(\Gamma_1, \forall x A, \Delta_1) \cap L(\Gamma_2, \Delta_2)$ , consider that by the induction hypothesis,  $L(I) \subseteq L(\Gamma_1, A[x/t], \Delta_1) \cap L(\Gamma_2, \Delta_2)$ . As there are no function symbols and since constant symbols are treated as function symbols,  $L(\forall x A) = L(A[x/t])$ .

2.  $\langle (\Gamma_1; \Delta_1), (\Gamma_2, \forall x A; \Delta_2) \rangle$ . This case can be argued analogously.

- The last rule application is an instance of  $\forall : r$ . Then it is of the form:

$$\frac{\Gamma \vdash \Delta, A[x/y]}{\Gamma \vdash \Delta, \forall x A} \forall : r$$

where  $y$  does not appear in  $\Gamma, \Delta$  or  $A$ .

There are two possible partition schemes of  $\Gamma \vdash \Delta, \forall x A$ :

1.  $\chi = \langle (\Gamma_1; \Delta_1, \forall x A), (\Gamma_2; \Delta_2) \rangle$ . By the induction hypothesis, there exists an interpolant  $I$  of the partition  $\langle (\Gamma_1; \Delta_1, A[x/y]), (\Gamma_2; \Delta_2) \rangle$  of the upper sequent. Hence for  $I$ ,  $\Gamma_1 \vdash \Delta_1, A[x/y], I$  and  $I, \Gamma_2 \vdash \Delta_2$  are provable. As  $y$  does not occur in  $\Gamma$  or  $\Delta$  and consequently by condition 3 does not occur in  $I$ , we may apply the  $\forall : r$  rule to the former sequent to obtain  $\Gamma_1 \vdash \Delta_1, \forall x A, I$ . Hence  $I$  is an interpolant for  $\chi$  as well.

2.  $\langle(\Gamma_1; \Delta_1), (\Gamma_2; \Delta_2, \forall x A)\rangle$ . This case can be argued analogously.  $\square$

This allows us to state the central theorem of this section:

**Theorem 3.11.** *Let  $\Gamma$  and  $\Delta$  be sets of first-order clauses without equality and function symbols such that  $\Gamma \cup \Delta$  is unsatisfiable. Then there is an interpolant for  $\Gamma$  and  $\Delta$ .*

*Proof.* We show that there is an interpolant for  $\Gamma \models \neg\Delta$ , which by Proposition 1.4 proves the theorem. By the completeness of cut-free sequent calculus, there is a proof of  $\Gamma \vdash \neg\Delta$ . By Lemma 3.10, there is an interpolant  $I$  for the partition  $\langle(\Gamma; ), (; \neg\Delta)\rangle$ .  $I$  is the desired interpolant for  $\Gamma \models \neg\Delta$ .  $\square$

### 3.3 Proof by reduction

Using the results of the previous sections, we can now give a proof of the interpolation theorem:

*Proof of Theorem 1.3 (Interpolation).* Since  $\Gamma \cup \Delta$  is unsatisfiable, by Proposition 3.7,  $T_{Ax}(\Gamma \cup \Delta)$  is unsatisfiable.

$$\begin{aligned} T_{Ax}(\Gamma \cup \Delta) &\Leftrightarrow \{F_{Ax}(L(\Gamma \cup \Delta)), E_{Ax}(L(\Gamma \cup \Delta))\} \cup T(\Gamma \cup \Delta) \\ &\Leftrightarrow \{F_{Ax}(L(\Gamma) \cup L(\Delta)), E_{Ax}(L(\Gamma) \cup L(\Delta))\} \cup T(\Gamma) \cup T(\Delta) \\ &\Leftrightarrow \{F_{Ax}(L(\Gamma)) \wedge F_{Ax}(L(\Delta)), E_{Ax}(L(\Gamma)) \wedge E_{Ax}(L(\Delta))\} \cup T(\Gamma) \cup T(\Delta) \\ &\Leftrightarrow \{F_{Ax}(L(\Gamma)), E_{Ax}(L(\Gamma))\} \cup T(\Gamma) \cup \{F_{Ax}(L(\Delta)), E_{Ax}(L(\Delta))\} \cup T(\Delta) \\ &\Leftrightarrow T_{Ax}(\Gamma) \cup T_{Ax}(\Delta) \end{aligned}$$

Hence  $T_{Ax}(\Gamma) \cup T_{Ax}(\Delta)$  is unsatisfiable as well. By Lemma 3.6.1  $T_{Ax}(\Gamma)$  and  $T_{Ax}(\Delta)$  contain neither function symbols nor the equality symbol. Hence by Theorem 3.11, there is an interpolant  $I$  such that

1.  $T_{Ax}(\Gamma) \models I$
2.  $T_{Ax}(\Delta) \models \neg I$
3.  $L(I) \subseteq L(T_{Ax}(\Gamma)) \cap L(T_{Ax}(\Delta))$

We now show that  $T^{-1}(I)$  is an interpolant for  $\Gamma$  and  $\Delta$ .

$T_{Ax}(\Gamma) \models I$  is equivalent to  $T_{Ax}(\Gamma) \cup \{\neg I\}$  being unsatisfiable. Through the unfolding of  $T_{Ax}(\Gamma)$ , we get that  $\{F_{Ax}(L(\Gamma)), E_{Ax}(L(\Gamma))\} \cup T(\Gamma) \cup \{\neg I\}$  is unsatisfiable. This set of formulas can now be translated back to the original language with the equality symbol and function symbols. More formally, since  $L(\neg I) \subseteq L(T_{Ax}(\Gamma))$ , we can apply Proposition 3.7.2 by considering  $T(\Gamma) \cup \{\neg I\}$  as  $\Phi$  to conclude that  $T^{-1}(T(\Gamma) \cup \{\neg I\})$  is unsatisfiable. By pulling  $T^{-1}$  inward and an application of Lemma 3.4, we get that  $\Gamma \cup \{T^{-1}(\neg I)\} = \Gamma \cup \{\neg T^{-1}(I)\}$  is unsatisfiable. Therefore  $\Gamma \models T^{-1}(I)$ .



For  $\Delta$ , an analogous argument goes through and so from  $T_{Ax}(\Gamma) \models \neg I$  we can deduce that  $\Delta \models \neg T^{-1}(I)$ .

By item 3,  $I$  is in the language  $L(T_{Ax}(\Gamma)) \cap L(T_{Ax}(\Delta))$ , which by Lemma 3.6.1 is  $T(L(\Gamma)) \cap T(L(\Delta))$ .

$$\begin{aligned}
 T(L(\Gamma)) \cap T(L(\Delta)) &= (L(\Gamma) \cup \{E\} \cup \{F_f \mid f \in FS(\Gamma)\}) \setminus (\{=\} \cup FS(\Gamma)) \cap \\
 &\quad (L(\Delta) \cup \{E\} \cup \{F_f \mid f \in FS(\Delta)\}) \setminus (\{=\} \cup FS(\Delta)) \\
 &= ((L(\Gamma) \cap L(\Delta)) \cup \{E\} \cup \{F_f \mid f \in FS(\Gamma) \cap FS(\Delta)\}) \setminus (\{=\} \cup FS(\Gamma) \cup FS(\Delta)) \\
 &= ((L(\Gamma) \cap L(\Delta)) \cup \{E\} \cup \{F_f \mid f \in FS(L(\Gamma) \cap L(\Delta))\}) \setminus (\{=\} \cup FS(L(\Gamma) \cap L(\Delta))) \\
 &= T(L(\Gamma) \cap L(\Delta))
 \end{aligned}$$

As  $I$  is in the language  $T(L(\Gamma) \cap L(\Delta))$ , by Lemma 3.6.2,  $T^{-1}(I)$  is in the language  $L(\Gamma) \cap L(\Delta)$ .  $\square$

# Proofs

## 4.1 WT: Interpolation extraction in one pass

easy for constants, just as in huang but in one pass

terms can grow unpredictably, order cannot be determined during pass

## 4.2 WT: Interpolation extraction in two passes

### 4.2.1 huang proof: propositional

Let  $\Gamma \cup \Delta$  be unsatisfiable. Let  $\pi$  be a proof of the empty clause from  $\Gamma \cup \Delta$ . Then PI is a function that returns a interpolant with respect to the current clause.

**Definition 4.1** (Propositional interpolant). Let  $\pi$  be a resolution refutation of  $\Gamma \cup \Delta$ . A formula  $A$  is a *propositional interpolant* if

1.  $\Gamma \models A$
2.  $\Delta \models \neg A$
3.  $\text{PS}(A) \subseteq (\text{PS}(\Gamma) \cap \text{PS}(\Delta)) \cup \{\top, \perp\}$ .

For a clause  $C$  in  $\pi$ , a formula  $A_C$  is a *propositional interpolant relative to  $C$*  if

1.  $\Gamma \models A_C \vee C$
2.  $\Delta \models \neg A_C \vee C$
3.  $\text{PS}(A_C) \subseteq (\text{PS}(\Gamma) \cap \text{PS}(\Delta)) \cup \{\top, \perp\}$ .

The propositional interpolant of the empty clause derived in  $\pi$  is denoted by  $\text{PI}(\pi)$ .

$\triangle$

The third condition of a propositional interpolant will sometimes be referred to as *language restriction*. It is easy to see that the propositional interpolant relative to the empty clause of a resolution refutation is a propositional interpolant.

We proceed by defining a procedure PI which extracts propositional interpolants from a resolution refutation.

**Definition 4.2** (Propositional interpolant extraction.). Let  $\pi$  be a resolution refutation of  $\Gamma \cup \Delta$ .  $\text{PI}(\pi)$  is defined to be  $\text{PI}(\square)$ , where  $\square$  is the empty clause derived in  $\pi$ .

For a clause  $C$  in  $\pi$ ,  $\text{PI}(C)$  is defined as follows:

Base case. If  $C \in \Gamma$ ,  $\text{PI}(C) = \perp$ . If otherwise  $C \in \Delta$ ,  $\text{PI}(C) = \top$ .

Resolution. If the clause  $C$  is the result of a resolution step of  $C_1 : D \vee l$  and  $C_2 : E \vee \neg l'$  using a unifier  $\sigma$  such that  $l\sigma = l'\sigma$ , then  $\text{PI}(C)$  is defined as follows:

1. If  $l$  is  $\Gamma$ -colored:  $\text{PI}(C) = [\text{PI}(C_1) \vee \text{PI}(C_2)]\sigma$
2. If  $l$  is  $\Delta$ -colored:  $\text{PI}(C) = [\text{PI}(C_1) \wedge \text{PI}(C_2)]\sigma$
3. If  $l$  is grey:  $\text{PI}(C) = [(l \wedge \text{PI}(C_2)) \vee (\neg l' \wedge \text{PI}(C_1))]\sigma$

Factorisation. If the clause  $C$  is the result of a factorisation of  $C_1 : l \vee l' \vee D$  using a unifier  $\sigma$  such that  $l\sigma = l'\sigma$ , then  $\text{PI}(C) = \text{PI}(C_1)\sigma$ .

Paramodulation. If the clause  $C$  is the result of a paramodulation of  $C_1 : s = t \vee C$  and  $C_2 : D[r]$  using a unifier  $\sigma$  such that  $r\sigma = s\sigma$ , then  $\text{PI}(C)$  is defined according to the following case distinction:

1. If  $r$  occurs in a maximal  $\Delta$ -term  $h(r)$  in  $D[r]$  and  $h(r)$  occurs more than once in  $D[r] \vee \text{PI}(D[r])$ :  
 $\text{PI}(C) = [(s = t \wedge \text{PI}(C_2)) \vee (s \neq t \wedge \text{PI}(C_1))]\sigma \vee (s = t \wedge h[s] \neq h[t])\sigma$
2. If  $r$  occurs in a maximal  $\Gamma$ -term  $h(r)$  in  $D[r]$  and  $h(r)$  occurs more than once in  $D[r] \vee \text{PI}(D[r])$ :  
 $\text{PI}(C) = [(s = t \wedge \text{PI}(C_2)) \vee (s \neq t \wedge \text{PI}(C_1))]\sigma \wedge (s \neq t \vee h[s] = h[t])\sigma$
3. Otherwise:  
 $\text{PI}(C) = [(s = t \wedge \text{PI}(C_2)) \vee (s \neq t \wedge \text{PI}(C_1))]\sigma \quad \triangle$

**Proposition 4.3.** *Let  $C$  be a clause of a resolution refutation of  $\Gamma \cup \Delta$ . Then  $\text{PI}(C)$  is a propositional interpolant with respect to  $C$ .*

*Proof.* Proof by induction on the number of rule applications including the following strengthenings:  $\Gamma \models \text{PI}(C) \vee C_\Gamma$  and  $\Delta \models \neg \text{PI}(C) \vee C_\Delta$ , where  $D_\Phi$  denotes the clause  $D$  with only the literals which are contained in  $L(\Phi)$ . They clearly imply conditions 1 and 2 of definition 4.1.

Base case. Suppose no rules were applied. We distinguish two possible cases:

1.  $C \in \Gamma$ . Then  $\text{PI}(C) = \perp$ . Clearly  $\Gamma \models \perp \vee C_\Gamma$  as  $C_\Gamma = C \in \Gamma$ ,  $\Delta \models \neg \perp \vee C_\Delta$  and  $\perp$  satisfies the restriction on the language.
2.  $C \in \Delta$ . Then  $\text{PI}(C) = \top$ . Clearly  $\Gamma \models \top \vee C_\Gamma$ ,  $\Delta \models \neg \top \vee C_\Delta$  as  $C_\Delta = C \in \Delta$  and  $\top$  satisfies the restriction on the language.

Suppose the property holds for  $n$  rule applications. We show that it holds for  $n+1$  applications by considering the last one:

Resolution. Suppose the last rule application is an instance of resolution. Then it is of the form:

$$\frac{C_1 : D \vee l \quad C_2 : E \vee \neg l'}{C : (D \vee E)\sigma} \quad l\sigma = l'\sigma$$

By the induction hypothesis, we can assume that:

$$\begin{aligned} \Gamma &\models \text{PI}(C_1) \vee (D \vee l)_\Gamma \\ \Delta &\models \neg \text{PI}(C_1) \vee (D \vee l)_\Delta \\ \Gamma &\models \text{PI}(C_2) \vee (E \vee \neg l')_\Gamma \\ \Delta &\models \neg \text{PI}(C_2) \vee (E \vee \neg l')_\Delta \end{aligned}$$

We consider the respective cases from definition 4.2:

1.  $l$  is  $\Gamma$ -colored. Then  $\text{PI}(C) = [\text{PI}(C_1) \vee \text{PI}(C_2)]\sigma$ .  
 As  $\text{PS}(l) \in L(\Gamma)$ ,  $\Gamma \models (\text{PI}(C_1) \vee D_\Gamma \vee l)\sigma$  as well as  $\Gamma \models (\text{PI}(C_2) \vee E_\Gamma \vee \neg l')\sigma$ .  
 By a resolution step, we get  $\Gamma \models (\text{PI}(C_1) \vee \text{PI}(C_2))\sigma \vee ((D \vee E)\sigma)_\Gamma$ .  
 Furthermore, as  $\text{PS}(l) \notin L(\text{PI})$ ,  $\Delta \models (\neg \text{PI}(C_1) \vee D_\Delta)\sigma$  as well as  $\Delta \models (\neg \text{PI}(C_2) \vee E_\Delta)\sigma$ . Hence it certainly holds that  $\Delta \models (\neg \text{PI}(C_1) \vee \neg \text{PI}(C_2))\sigma \vee (D \vee E)\sigma_\Delta$ .  
 The language restriction clearly remains satisfied as no non-logical symbols are added.
2.  $l$  is  $\Delta$ -colored. Then  $\text{PI}(C) = [\text{PI}(C_1) \wedge \text{PI}(C_2)]\sigma$ .  
 As  $\text{PS}(l) \notin L(\Gamma)$ ,  $\Gamma \models (\text{PI}(C_1) \vee D_\Gamma)\sigma$  as well as  $\Gamma \models (\text{PI}(C_2) \vee E_\Gamma)\sigma$ . Suppose that in a model  $M$  of  $\Gamma$ ,  $M \not\models D_\Gamma$  and  $M \not\models E_\Gamma$ . Then  $M \models \text{PI}(C_1) \wedge \text{PI}(C_2)$ . Hence  $\Gamma \models (\text{PI}(C_1) \wedge \text{PI}(C_2))\sigma \vee ((D \vee E)\sigma)_\Gamma$ .  
 Furthermore due to  $\text{PS}(l) \in L(\Delta)$ ,  $\Delta \models (\neg \text{PI}(C_1) \vee D_\Delta \vee l)\sigma$  as well as  $\Delta \models (\neg \text{PI}(C_2) \vee E_\Delta \vee \neg l')\sigma$ . By a resolution step, we get  $\Delta \models (\neg \text{PI}(C_1) \vee \neg \text{PI}(C_2))\sigma \vee (D_\Delta \vee E_\Delta)\sigma$  and hence  $\Delta \models \neg(\text{PI}(C_1) \wedge \text{PI}(C_2))\sigma \vee (D_\Delta \vee E_\Delta)\sigma$ .  
 The language restriction again remains intact.
3.  $l$  is grey. Then  $\text{PI}(C) = [(l \wedge \text{PI}(C_2)) \vee (\neg l' \wedge \text{PI}(C_1))]\sigma$ .  
 First, we have to show that  $\Gamma \models [(l \wedge \text{PI}(C_2)) \vee (l' \wedge \text{PI}(C_1))]\sigma \vee ((D \vee E)\sigma)_\Gamma$ .  
 Suppose that in a model  $M$  of  $\Gamma$ ,  $M \not\models D_\Gamma$  and  $\Gamma \not\models E$ . Otherwise we are done.  
 The induction assumption hence simplifies to  $M \models \text{PI}(C_1) \vee l$  and  $M \models$

$\text{PI}(C_2) \vee \neg l'$  respectively. As  $l\sigma = l'\sigma$ , by a case distinction argument on the truth value of  $l\sigma$ , we get that either  $M \models (l \wedge \text{PI}(C_2))\sigma$  or  $M \models (\neg l' \wedge \text{PI}(C_1))\sigma$ . Second, we show that  $\Delta \models ((l \vee \neg \text{PI}(C_1)) \wedge (\neg l' \vee \neg \text{PI}(C_2)))\sigma \vee ((D \vee E)\sigma)_\Delta$ . Suppose again that in a model  $M$  of  $\Delta$ ,  $M \not\models D_\Delta$  and  $\Gamma \not\models E_\Delta$ . Then the required statement follows from the induction hypothesis.

The language condition remains satisfied as only the common literal  $l$  is added to the interpolant.

**Factorisation.** Suppose the last rule application is an instance of factorisation. Then it is of the form:

$$\frac{C_1 : l \vee l' \vee D}{C_1 : (l \vee D)\sigma} \quad \sigma = \text{mgu}(l, l')$$

Then the propositional interpolant  $\text{PI}(C)$  is defined as  $\text{PI}(C_1)$ . By the induction hypothesis, we have:

$$\Gamma \models \text{PI}(C_1) \vee (l \vee l' \vee D)_\Gamma$$

$$\Delta \models \text{PI}(C_1) \vee (l \vee l' \vee D)_\Delta$$

It is easy to see that then also:

$$\Gamma \models (\text{PI}(C_1) \vee (l \vee D)_\Gamma)\sigma$$

$$\Delta \models (\text{PI}(C_1)\sigma \vee (l \vee D)_\Delta)\sigma$$

The restriction on the language trivially remains intact.

**Paramodulation.** Suppose the last rule application is an instance of paramodulation. Then it is of the form:

$$\frac{C_1 : D \vee s = t \quad C_2 : E[r]}{C : (D \vee E[t])\sigma} \quad \sigma = \text{mgu}(s, r)$$

By the induction hypothesis, we have:

$$\Gamma \models \text{PI}(C_1) \vee (D \vee s = t)_\Gamma$$

$$\Delta \models \neg \text{PI}(C_1) \vee (D \vee s = t)_\Delta$$

$$\Gamma \models \text{PI}(C_2) \vee (E[r])_\Gamma$$

$$\Delta \models \neg \text{PI}(C_2) \vee (E[r])_\Delta$$

First, we show that  $\text{PI}(C)$  as constructed in case 3 of the definition is a propositional interpolant in any of these cases:

$$\text{PI}(C) = (s = t \wedge \text{PI}(C_2)) \vee (s \neq t \wedge \text{PI}(C_1))$$

Suppose that in a model  $M$  of  $\Gamma$ ,  $M \not\models D\sigma$  and  $M \not\models E[t]\sigma$ . Otherwise we are done. Furthermore, assume that  $M \models (s = t)\sigma$ . Then  $M \not\models E[r]\sigma$ , but then necessarily  $M \models \text{PI}(C_2)\sigma$ .

On the other hand, suppose  $M \models (s \neq t)\sigma$ . As also  $M \not\models D\sigma$ ,  $M \models \text{PI}(C_1)\sigma$ . Consequently,  $M \models [(s = t \wedge \text{PI}(C_2)) \vee (s \neq t \wedge \text{PI}(C_1))]\sigma \vee [(D \vee E)_\Gamma]\sigma$

By an analogous argument, we get  $\Delta \models [(s = t \wedge \neg \text{PI}(C_2)) \vee (s \neq t \wedge \neg \text{PI}(C_1))]\sigma \vee [(D \vee E)_\Delta]\sigma$ , which implies  $\Delta \models [(s \neq t \vee \neg \text{PI}(C_2)) \wedge (s = t \vee \neg \text{PI}(C_1))]\sigma \vee ((D \vee E)_\Delta)\sigma$

The language restriction again remains satisfied as the only predicate, that is added to the interpolant, is  $=$ .

This concludes the argumentation for case 3.

The interpolant for case 1 differs only by an additional formula added via a disjunction and hence condition 1 of definition 4.1 holds by the above reasoning. As the adjoined formula is a contradiction, its negation is valid which in combination with the above reasoning establishes condition 2. Since no new predicates are added, the language condition remains intact.

The situation in case 2 is somewhat symmetric: As a tautology is added to the interpolant with respect to case 1, condition 1 is satisfied by the above reasoning. For condition 2, consider that the negated interpolant for case 1 implies the negated interpolant for this case. The language condition again remains intact.  $\square$

#### 4.2.2 huang proof: overbinding

Before we are able to specify a procedure to transform the propositional interpolant generated by PI into a proper interpolant without any colored terms, we define a simpler but equally powerful form of resolution refutations.

**Definition 4.4.** A resolution refutation is a *tree refutation* if every clause is used at most once.  $\triangle$

The following lemma shows that this form does not restrict the calculus given that we allow multisets as initial clause sets.

**Lemma 4.5.** *Every resolution refutation can be transformed into a tree refutation.*

*Proof.* Let  $\pi$  be a resolution refutation of  $\Phi$ . We prove that  $\pi$  can be transformed into a tree refutation by induction on the number of clauses that are used multiple times.

Suppose that no clause is used more than once in  $\pi$ . Then  $\pi$  is a tree refutation.

Otherwise let  $\Psi$  be the set of clauses which is used multiple times. Let  $C \in \Psi$  be such that no clause  $D \in \Psi$  is used in the derivation leading to  $C$ . Let  $\chi$  be the derivation leading to  $C$ .

Suppose  $C$  is used  $m$  times. We create another resolution refutation  $\pi'$  from  $\pi$  which contains  $m$  copies of  $\chi$  and replaces the  $i$ th use of the clause  $C$  by the final clause of the  $i$ th copy of  $\chi$ ,  $1 \leq i \leq m$ . In order to ensure that the sets of variables of the input clauses are disjoint, we rename the variables in each copy of  $\chi$  and adapt  $\pi'$  accordingly. Hence  $\pi'$  is a resolution refutation of  $\Phi$  where  $m - 1$  clauses are used more than once.  $\square$

In a tree refutation where the input clauses have a disjoint sets of variables, every variable has a unique ancestor which traces back to an input clause and hence appears only along a certain path. This insight allows us to push substitutions of the variables upwards along this path and arrive at the following definition and lemma:

**Definition 4.6.** A resolution refutation is a *propositional refutation* if no nontrivial substitutions are employed.  $\triangle$

**Lemma 4.7.** *Let  $\Phi$  be unsatisfiable. Then there is a propositional refutation of  $\Phi$  which starts from instances of  $\Phi$ .*

*Proof.* Let  $\pi$  be a resolution refutation of  $\Phi$ . By Lemma 4.5, we can assume without loss of generality that  $\pi$  is a tree refutation where the sets of variables of the input clauses are disjoint. Furthermore, we can assume that only most general unifiers are employed in  $\pi$ .

Then any unifier in  $\pi$  is either trivial on  $x$  or there is one unique unifier  $\sigma$  in  $\pi$  with  $x\sigma = t$  where  $x$  does not occur in  $t$ . Hence along the path through the deduction where  $x$  occurs, it remains unchanged. Therefore we can create a new resolution refutation  $\pi'$  from  $\pi$  where  $x$  is replaced by  $t$ . Clearly  $\pi'$  is rooted in instances of  $\Phi$ .

By application of this procedure to all variable occurring in  $\pi$ , we obtain a desired resolution refutation.  $\square$

Even though propositional refutations have nice properties for theoretical analysis, their use in practise is not desired as its construction involves a considerable blowup of the refutation. But its use is still justified in this instance as we can show for arbitrary refutations  $\pi$  that the algorithm stated in 4.2 gives closely related results for both  $\pi$  and its corresponding propositional refutation.

**Lemma 4.8.** *Let  $\pi$  be a resolution refutation of  $\Phi$  and  $\pi'$  a propositional refutation corresponding to  $\pi$ . Then for every clause  $C$  in  $\pi$  and its corresponding clause  $C'$  in  $\pi'$ ,  $\text{PI}(C)\sigma = \text{PI}(C')$ , where  $\sigma$  is the composition of the unifications of  $\pi$  which are applied to the variables occurring in  $C$ .*

*Proof.* For the construction of the propositional skeleton of  $\text{PI}(\cdot)$  only the coloring of the clauses is relevant and since this is the same in both  $\pi$  and  $\pi'$ , it coincides for  $\text{PI}(C)$  and  $\text{PI}(C')$ .

Hence  $\text{PI}(C)$  and  $\text{PI}(C')$  differ only in their term structure. To be more specific, in  $\text{PI}(C')$ , the composition of substitutions that are applied in  $\pi$  have already been applied to the initial clauses of  $\pi'$ . Note that substitution commutes with the rules of resolution. Therefore the only difference between  $\text{PI}(C)$  and  $\text{PI}(C')$  is that at certain term positions, there are variables in  $\text{PI}(C)$  where in  $\text{PI}(C')$  by some substitution a different term is located. But these substitutions are certainly applied by  $\sigma$ , hence  $\text{PI}(C)\sigma = \text{PI}(C')$ .  $\square$

This establishes the theoretical framework which is required to define and show the correctness of the procedure to construct a proper interpolant from the propositional interpolant. The idea of this procedure will be to replace colored terms still occurring

in the propositional interpolant with variables and quantifying them appropriately. This replacement is referred to as lifting:

**Definition 4.9** (Lifting). Let  $\Gamma$  and  $\Delta$  be sets of first-order formulas,  $\phi$  a formula or a term,  $t_1, \dots, t_n$  the maximal  $\Phi$ -terms for  $\Phi \in \{\Gamma, \Delta\}$  in  $\phi$  and  $x_1, \dots, x_n$  fresh variables. Then  $\ell_{\Phi, x}[\phi]$  denotes  $\phi\{t_1/x_1\} \dots \{t_n/x_n\}$ .  $\triangle$

**Lemma 4.10.** *Let  $A$  and  $B$  be first-order formulas. Then it holds that:*

1.  $\ell_{\Phi, x}[\neg A] \Leftrightarrow \neg \ell_{\Phi, x}[A]$
2.  $\ell_{\Phi, x}[A \circ B] \Leftrightarrow (\ell_{\Phi, x}[A] \circ \ell_{\Phi, x}[B])$  for  $\circ \in \{\wedge, \vee\}$

First, we consider the lifting of the  $\Delta$ -terms:

**Lemma 4.11.** *Let  $\pi$  be a resolution refutation of  $\Gamma \cup \Delta$ . Then  $\Gamma \models \ell_{\Delta, x}[\text{PI}(C) \vee C]$  for  $C$  in  $\pi$ .*

*Proof.* We proof this result by induction on the number of rule applications in the propositional refutation corresponding to  $\pi$ . Similar to the proof of 4.3, we show the strengthening:  $\Gamma \models \ell_{\Delta, x}[\text{PI}(C) \vee C_\Gamma]$  for  $C$  in  $\pi$ .

Base case. If no rules have been applied,  $C$  is an instance of a clause of either  $\Gamma$  or  $\Delta$ .

In the former case, all  $\Delta$ -terms of  $C$  were added by unification, hence by replacing them with variables, we obtain a clause  $C'$  which still is an instance of  $C$  and consequently is implied by  $\Gamma$ . In the latter case,  $\text{PI}(C) = \top$ .

Resolution. Suppose the last rule application is an instance of resolution. Then it is of the form:

$$\frac{C_1 : D \vee l \quad C_2 : E \vee \neg l}{C : D \vee E}$$

By the induction hypothesis,

$$\Gamma \models \ell_{\Delta, x}[\text{PI}(C_1) \vee (D \vee l)_\Gamma] \text{ and}$$

$$\Gamma \models \ell_{\Delta, x}[\text{PI}(C_2) \vee (E \vee \neg l)_\Gamma]$$

which by Lemma 4.10 is equivalent to

$$\Gamma \models \ell_{\Delta, x}[\text{PI}(C_1)] \vee \ell_{\Delta, x}[D_\Gamma] \vee \ell_{\Delta, x}[l_\Gamma] \quad (^{\circ}) \text{ and}$$

$$\Gamma \models \ell_{\Delta, x}[\text{PI}(C_2)] \vee \ell_{\Delta, x}[E_\Gamma] \vee \neg \ell_{\Delta, x}[l_\Gamma] \quad (^*) .$$

1. Suppose  $l$  is  $\Gamma$ -colored. Then  $\text{PI}(C) = \text{PI}(C_1) \vee \text{PI}(C_2)$ . By using resolution of  $(^*)$  and  $(^{\circ})$  on  $\ell_{\Delta, x}[l_\Gamma]$ , we get that

$$\Gamma \models \ell_{\Delta, x}[\text{PI}(C_1)] \vee \ell_{\Delta, x}[\text{PI}(C_2)] \vee \ell_{\Delta, x}[D_\Gamma] \vee \ell_{\Delta, x}[E_\Gamma].$$

Several applications of Lemma 4.10 give  $\Gamma \models \ell_{\Delta, x}[\text{PI}(C_1) \vee \text{PI}(C_2) \vee (D \vee E)_\Gamma]$ .



2. Suppose  $l$  is  $\Delta$ -colored. Then  $\text{PI}(C) = \text{PI}(C_1) \wedge \text{PI}(C_2)$ .

As  $l$  and  $\neg l$  are not contained in  $L(\Gamma)$ , we get that

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_1)] \vee \ell_{\Delta,x}[D_\Gamma] \text{ and}$$

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_2)] \vee \ell_{\Delta,x}[E_\Gamma].$$

So if in a model  $M$  of  $\Gamma$  we have that  $M \not\models \ell_{\Delta,x}[D_\Gamma]$  and  $M \not\models \ell_{\Delta,x}[E_\Gamma]$ , it follows that  $M \models \ell_{\Delta,x}[\text{PI}(C_1)]$  and  $M \models \ell_{\Delta,x}[\text{PI}(C_2)]$ . Hence by Lemma 4.10  $M \models \ell_{\Delta,x}[\text{PI}(C_1) \wedge \text{PI}(C_2)] \vee \ell_{\Delta,x}[(D \vee E)_\Gamma]$ .

3. Suppose  $l$  is grey. Then  $\text{PI}(C) = (l \wedge \text{PI}(C_2)) \vee (\neg l \wedge \text{PI}(C_1))$ .

We show that  $\Gamma \models \ell_{\Delta,x}[(l \wedge \text{PI}(C_2)) \vee (\neg l \wedge \text{PI}(C_1)) \vee (D \vee E)_\Gamma]$ .

Suppose that for a model  $M$  of  $\Gamma$  that  $M \not\models \ell_{\Delta,x}[D_\Gamma]$  and  $M \not\models \ell_{\Delta,x}[E_\Gamma]$ . Then by  $(\circ)$  and  $(*)$ , we get that

$$M \models \ell_{\Delta,x}[\text{PI}(C_1)] \vee \ell_{\Delta,x}[l_\Gamma] \text{ as well as}$$

$$M \models \ell_{\Delta,x}[\text{PI}(C_2)] \vee \neg \ell_{\Delta,x}[l_\Gamma].$$

So  $M \models \ell_{\Delta,x}[l_\Gamma]$  implies that  $M \models \ell_{\Delta,x}[\text{PI}(C_2)]$  and  $M \models \neg \ell_{\Delta,x}[l_\Gamma]$  implies that  $M \models \ell_{\Delta,x}[\text{PI}(C_1)]$  and

Therefore  $M \models (\ell_{\Delta,x}[l] \wedge \ell_{\Delta,x}[\text{PI}(C_2)]) \vee (\neg \ell_{\Delta,x}[l] \wedge \ell_{\Delta,x}[\text{PI}(C_1)]) \vee (\ell_{\Delta,x}[D_\Gamma] \vee \ell_{\Delta,x}[E_\Gamma])$ , and several applications of Lemma 4.10 give  $M \models \ell_{\Delta,x}[(l \wedge \text{PI}(C_2)) \vee (\neg l \wedge \text{PI}(C_1)) \vee (D_\Gamma \vee E_\Gamma)]$ .

Factorisation. Suppose the last rule application is an instance of factorisation. Then it is of the form:

$$\frac{C_1 : l \vee l \vee D}{C : l \vee D}$$

The propositional interpolant directly carried over from  $C_1$ , i.e.  $\text{PI}(C) = \text{PI}(C_1)$ .

By the induction hypothesis, we get that  $\Gamma \models \ell_{\Delta,x}[\text{PI}(C_1) \vee (l \vee l \vee D)_\Gamma]$ . By Lemma 4.10,

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_1)] \vee (\ell_{\Delta,x}[l_\Gamma] \vee \ell_{\Delta,x}[l_\Gamma] \vee \ell_{\Delta,x}[D_\Gamma]),$$

which clearly is equivalent to

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_1)] \vee (\ell_{\Delta,x}[l_\Gamma] \vee \ell_{\Delta,x}[D_\Gamma]),$$

so by again applying Lemma 4.10, we arrive at

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_1) \vee (l \vee D)_\Gamma].$$

Paramodulation. Suppose the last rule application is an instance of paramodulation.

Then it is of the form:

$$\frac{C_1 : D \vee s = t \quad C_2 : E[s]_p}{C : D \vee E[t]_p}$$

By the induction hypothesis, we have that

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_1) \vee (D \vee s = t)_\Gamma] \text{ and}$$

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_2) \vee (E[s]_p)_\Gamma].$$

By Lemma 4.10, we get that

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_1)] \vee \ell_{\Delta,x}[D_\Gamma] \vee \ell_{\Delta,x}[s] = \ell_{\Delta,x}[t] \text{ and}$$

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_2)] \vee \ell_{\Delta,x}[(E[s]_p)_\Gamma].$$

We distinguish two cases:

1. Suppose  $s$  does not occur in a maximal  $\Delta$ -term  $h[s]$  in  $E[s]_p$  which occurs more than once in  $\text{PI}(E(s)) \vee E[s]_p$ .

We show that  $\Gamma \models \ell_{\Delta,x}[(s = t \wedge \text{PI}(C_2)) \vee (s \neq t \wedge \text{PI}(C_1)) \vee (D \vee E[t]_p)_\Gamma]$ , which subsumes the cases 2 and 3 of Definition 4.2. By Lemma 4.10, this is equivalent to

$$\Gamma \models (\ell_{\Delta,x}[s] = \ell_{\Delta,x}[t] \wedge \ell_{\Delta,x}[\text{PI}(C_2)]) \vee (\ell_{\Delta,x}[s] \neq \ell_{\Delta,x}[t] \wedge \ell_{\Delta,x}[\text{PI}(C_1)]) \vee (\ell_{\Delta,x}[D_\Gamma] \vee \ell_{\Delta,x}[(E[t]_p)_\Gamma])$$

Suppose that in a model  $M$  of  $\Gamma$ ,  $M \not\models \ell_{\Delta,x}[D_\Gamma]$  and  $M \not\models \ell_{\Delta,x}[(E[t]_p)_\Gamma]$ . We show that then, depending on whether  $\ell_{\Delta,x}[s] = \ell_{\Delta,x}[t]$  holds in  $M$ , one of the first two disjuncts holds in  $M$ .

Then in case  $M \models \ell_{\Delta,x}[s] = \ell_{\Delta,x}[t]$  we also get  $M \not\models \ell_{\Delta,x}[(E[s]_p)_\Gamma]$  and consequently by the induction hypothesis  $M \models \ell_{\Delta,x}[\text{PI}(C_2)]$ .

However in case  $M \models \ell_{\Delta,x}[s] \neq \ell_{\Delta,x}[t]$  we get by the induction hypothesis that  $M \models \ell_{\Delta,x}[\text{PI}(C_1)]$ .

2. Otherwise  $s$  occurs in a maximal  $\Delta$ -term  $h[s]$  in  $E[s]_p$  which occurs more than once in  $\text{PI}(E(s)) \vee E[s]_p$ . This reflects case 1 of Definition 4.2.

Then models are possible in which  $s = t$  and therefore  $\ell_{\Delta,x}[s] = \ell_{\Delta,x}[t]$  holds, while at the same time  $\ell_{\Delta,x}[h[s]] \neq \ell_{\Delta,x}[h[t]]$  does not as  $h[s]$  and  $h[t]$  are replaced by distinct variables due to being different  $\Delta$ -terms.

Therefore we amend the proof of case 1 as follows:

In case  $M \models \ell_{\Delta,x}[s] = \ell_{\Delta,x}[t]$  (otherwise proceed as in case 1), one of the following cases holds:

- $M \models \ell_{\Delta,x}[h[s]] = \ell_{\Delta,x}[h[t]]$ . From this, it follows that as in the proof of case 1,  $M \not\models \ell_{\Delta,x}[(E[s]_p)_\Gamma]$  and consequently  $M \models \ell_{\Delta,x}[\text{PI}(C_2)]$  again by the induction hypothesis.
- $M \models \ell_{\Delta,x}[h[s]] \neq \ell_{\Delta,x}[h[t]]$ . However as here  $\text{PI}(C)$  contains the with respect to case 1 additional disjunct  $s = t \wedge h[s] \neq h[t]$ ,  $M \models \ell_{\Delta,x}[\text{PI}(C)]$  due to  $M \models \ell_{\Delta,x}[s] = \ell_{\Delta,x}[t] \wedge \ell_{\Delta,x}[h[s]] \neq \ell_{\Delta,x}[h[t]]$   $\square$

*Remark.* DOES NOT JUST WORK LIKE THIS

This lemma does not directly hold for a non-propositional proof.

The statement of the lemma would be identical in this case. For resolution, by the induction hypothesis, we would get

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_1)] \vee \ell_{\Delta,x}[D] \vee \ell_{\Delta,x}[l]$$

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_2)] \vee \ell_{\Delta,x}[E] \vee \neg \ell_{\Delta,x}[l']$$

By the resolution we know that  $l\sigma = l'\sigma$ .

In order to proceed as in the other proof, we would need  $\sigma'$  s.t.  $\ell_{\Delta,x}[l]\sigma' = \ell_{\Delta,x}[l']\sigma'$

$$\Gamma \models (\ell_{\Delta,x}[\text{PI}(C_1)] \vee \ell_{\Delta,x}[D] \vee \ell_{\Delta,x}[l])\sigma'$$

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_1)]\sigma' \vee \ell_{\Delta,x}[D]\sigma' \vee \ell_{\Delta,x}[l]\sigma'$$

$$\Gamma \models (\ell_{\Delta,x}[\text{PI}(C_2)] \vee \ell_{\Delta,x}[E] \vee \neg \ell_{\Delta,x}[l'])\sigma'$$

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_2)]\sigma' \vee \ell_{\Delta,x}[E]\sigma' \vee \neg \ell_{\Delta,x}[l']\sigma'$$

Hence

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_2)]\sigma' \vee \ell_{\Delta,x}[E]\sigma' \vee \ell_{\Delta,x}[\text{PI}(C_1)]\sigma' \vee \ell_{\Delta,x}[D]\sigma'$$

and therefore

$$\Gamma \models \ell_{\Delta,x}[\text{PI}(C_2) \vee E \vee \text{PI}(C_1) \vee D]\sigma'$$

And we need to show that  $\Gamma \models \ell_{\Delta,x}[(\text{PI}(C_1) \vee \text{PI}(C_2) \vee D \vee E)\sigma]$

### 4.2.3 final step of huang's proof

The definition PI possesses a convenient property which is termed *symmetry* in [DKPW10, Definition 5] and can be stated formally as follows:

**Lemma 4.12.** *Let  $\pi$  be a resolution refutation of  $\Gamma \cup \Delta$  and  $\hat{\pi}$  be  $\pi$  with  $\hat{\Gamma} = \Delta$  and  $\hat{\Delta} = \Gamma$ . Then  $\text{PI}(\pi) \Leftrightarrow \neg \text{PI}(\hat{\pi})$ .*

*Proof.* We prove this lemma by induction on  $\pi$ . Let  $\hat{\varphi}$  denote the clause/formula/literal/term in  $\hat{\pi}$  corresponding to the clause/formula/literal/term  $\varphi$  in  $\pi$ .

Base case. If  $C \in \Gamma$ , then  $C' \in \Delta'$  and  $\text{PI}(C) = \perp \Leftrightarrow \neg \top = \neg \text{PI}(C')$ . The case for  $C \in \Delta$  is analogous.

Resolution. If the clause  $C$  is the result of a resolution step of  $C_1 : D \vee l$  and  $C_2 : E \vee \neg l'$  using a unifier  $\sigma$  such that  $l\sigma = l'\sigma$  <sup>(\*)</sup>, then by induction hypothesis, we get that  $\text{PI}(C_i) = \neg \text{PI}(C'_i)$  for  $i \in \{1, 2\}$ .

We distinguish the following cases:

1.  $l$  is  $\Gamma$ -colored. Then  $\hat{l}$  is  $\Delta$ -colored.

$$\begin{aligned} \text{PI}(C) &= \text{PI}(C_1) \vee \text{PI}(C_2) \\ &\Leftrightarrow \neg(\neg \text{PI}(C_1) \wedge \neg \text{PI}(C_2)) \\ &= \neg(\text{PI}(\hat{C}_1) \wedge \text{PI}(\hat{C}_2)) \\ &= \neg \text{PI}(\hat{C}) \end{aligned}$$

2.  $l$  is  $\Delta$ -colored. This case can be argued analogously.

3.  $l$  is grey. Then  $\hat{l}$  is grey.

$$\begin{aligned}
\text{PI}(C) &= [(l \wedge \text{PI}(C_2)) \vee (\neg l' \wedge \text{PI}(C_1))] \sigma \\
&= (l\sigma \wedge \text{PI}(C_2)\sigma) \vee (\neg l'\sigma \wedge \text{PI}(C_1)\sigma) \\
&\Leftrightarrow (\neg l\sigma \vee \text{PI}(C_2)\sigma) \wedge (l'\sigma \vee \text{PI}(C_1)\sigma) \\
&\Leftrightarrow \neg[(l\sigma \wedge \neg \text{PI}(C_2)\sigma) \vee (\neg l'\sigma \wedge \neg \text{PI}(C_1)\sigma)] \\
&= \neg[(\hat{l}\sigma \wedge \neg \text{PI}(\hat{C}_2)\sigma) \vee (\neg \hat{l}'\sigma \wedge \neg \text{PI}(\hat{C}_1)\sigma)] \\
&= \neg[(\hat{l} \wedge \neg \text{PI}(\hat{C}_2)) \vee (\neg \hat{l}' \wedge \neg \text{PI}(\hat{C}_1))] \sigma \\
&= \neg[(\hat{l} \wedge \text{PI}(\hat{C}_2)) \vee (\neg \hat{l}' \wedge \text{PI}(\hat{C}_1))] \sigma \\
&= \neg \text{PI}(\hat{C})
\end{aligned}$$

Factorisation. Suppose the clause  $C$  is the result of a factorisation of  $C_1 : l \vee l' \vee D$ . Then  $\text{PI}(C) = \text{PI}(C_1)\sigma$  and the induction hypothesis gives the result.

Paramodulation. Suppose the clause  $C$  is the result of a paramodulation of  $C_1 : s = t \vee C$  and  $C_2 : D[r]$  using a unifier  $\sigma$  such that  $r\sigma = s\sigma$ . We distinguish the following cases:

1.  $r$  occurs in a maximal  $\Delta$ -term  $h[r]$  in  $D[r]$  and  $h[r]$  occurs more than once in  $D[r] \vee \text{PI}(D[r])$ . Then  $\hat{r}$  occurs in a maximal  $\Gamma$ -term  $\hat{h}[r]$  in  $\hat{D}[r]$  and  $\hat{h}[r]$  occurs more than once in  $\hat{D}[r] \vee \text{PI}(\hat{D}[r])$ .

$$\begin{aligned}
\text{PI}(C) &= [(s = t \wedge \text{PI}(C_2)) \vee (s \neq t \wedge \text{PI}(C_1))] \sigma \vee (s = t \wedge h[s] \neq h[t]) \sigma \\
&= [(s = t \wedge \neg \text{PI}(\hat{C}_2)) \vee (s \neq t \wedge \neg \text{PI}(\hat{C}_1))] \sigma \vee (s = t \wedge h[s] \neq h[t]) \sigma \\
&\Leftrightarrow \neg[(s \neq t \vee \text{PI}(\hat{C}_2)) \wedge (s = t \vee \text{PI}(\hat{C}_1))] \sigma \wedge \neg(s \neq t \vee h[s] = h[t]) \sigma \\
&\Leftrightarrow \neg[(s = t \wedge \text{PI}(\hat{C}_2)) \vee (s \neq t \wedge \text{PI}(\hat{C}_1))] \sigma \wedge \neg(s \neq t \vee h[s] = h[t]) \sigma \\
&= \neg \text{PI}(\hat{C})
\end{aligned}$$

2.  $r$  occurs in a maximal  $\Gamma$ -term  $h(r)$  in  $D[r]$  and  $h(r)$  occurs more than once in  $D[r] \vee \text{PI}(D[r])$ . This case can be argued analogously.
3. Otherwise:

$$\begin{aligned}
\text{PI}(C) &= [(s = t \wedge \text{PI}(C_2)) \vee (s \neq t \wedge \text{PI}(C_1))] \sigma \\
&= [(s = t \wedge \neg \text{PI}(\hat{C}_2)) \vee (s \neq t \wedge \neg \text{PI}(\hat{C}_1))] \sigma \\
&\Leftrightarrow \neg[(s \neq t \vee \text{PI}(\hat{C}_2)) \wedge (s = t \vee \text{PI}(\hat{C}_1))] \sigma \\
&\Leftrightarrow \neg[(s = t \wedge \text{PI}(\hat{C}_2)) \vee (s \neq t \wedge \text{PI}(\hat{C}_1))] \sigma \\
&= \neg \text{PI}(\hat{C})
\end{aligned}$$

□

This lemma can be leveraged to show a counterpart of Lemma 4.11 for  $\Delta$ :

**Corollary 4.13.** *Let  $\pi$  be a resolution refutation of  $\Gamma \cup \Delta$ . Then  $\Delta \models \ell_{\Gamma,x}[\neg \text{PI}(C) \vee C]$  for  $C$  in  $\pi$ .*

*Proof.* Build  $\hat{\pi}$  from  $\pi$  using  $\hat{\Gamma} = \Delta$  and  $\hat{\Delta} = \Gamma$  as initial clause set partition.

By Lemma 4.11,  $\hat{\Gamma} \models \ell_{\hat{\Delta},x}[\text{PI}(\hat{C}) \vee \hat{C}]$  for  $\hat{C}$  in  $\hat{\pi}$ .

By Lemma 4.12,  $\hat{\Gamma} \models \ell_{\hat{\Delta},x}[\neg \text{PI}(C) \vee \hat{C}]$  for the clause  $C$  in  $\pi$  corresponding to  $\hat{C}$  in  $\hat{\pi}$ . This however is nothing else than  $\Delta \models \ell_{\Gamma,x}[\neg \text{PI}(C) \vee C]$ .  $\square$

**Lemma 4.14.**  *$\ell_{\Gamma,y}[\ell_{\Delta,x}[C]]$  and  $\ell_{\Delta,x'}[\ell_{\Gamma,y'}[C]]$  differ only in the naming of the variables replacing maximal colored terms.*

*Proof.* Suppose a term  $t$  in  $C$  is affected by a lifting. We only need to consider maximal colored terms as grey terms are not affected by the liftings. Without loss of generality let  $t$  be a maximal  $\Delta$ -colored term.

Let  $\Phi$  be the positions of maximal occurrences of  $t$ . Then in the left hand side, exactly all terms at positions  $\Phi$  are replaced by  $x_i$  for some  $i$ .

In the right hand side, all terms at positions  $\Phi$  are replaced by  $\ell_{\Gamma,y'}[t]$  first. However after this step, all these terms are equal to  $\ell_{\Gamma,y'}[t]$ , and as all distinct maximal  $\Gamma$ -terms are replaced by distinct variables, no other maximal colored term is equal to  $\ell_{\Gamma,y'}[t]$ . Hence exactly the terms at positions  $\Phi$  are replaced by the same variable  $x'_j$  for some  $j$ .  $\square$

**Theorem 4.15.** *Let  $\pi$  be a resolution refutation of  $\Gamma \cup \Delta$  and  $z_1, \dots, z_n$  be the variables which replace the colored terms in  $\ell_{\Gamma,y}[\ell_{\Delta,x}[\text{PI}(\pi)]]$  ordered by their length. Then  $Q_1 z_1 \dots Q_n z_n \ell_{\Gamma,y}[\ell_{\Delta,x}[\text{PI}(\pi)]]$ , where  $Q_i$  is  $\forall$  ( $\exists$ ) if  $z_i$  replaces a  $\Delta$  ( $\Gamma$ )-term, is an interpolant.*

*Proof.* By Lemma 4.11,  $\Gamma \models \forall x_1 \dots \forall x_m \ell_{\Delta,x}[\text{PI}(\pi)]$  where  $m$  is the number of maximal  $\Delta$ -colored terms in  $\text{PI}(\pi)$ .

A term in  $\ell_{\Delta,x}[\text{PI}(\pi)]$  is either  $x_i$ ,  $1 \leq i \leq m$ , a grey term or a  $\Gamma$ -terms. Let  $t$  be a maximal  $\Gamma$ -term in  $\ell_{\Delta,x}[\text{PI}(\pi)]$  and  $x_{j_1}, \dots, x_{j_k}$  the variables replacing  $\Delta$ -terms in  $t$ . Note that the  $\Delta$ -terms, which are replaced by  $x_{j_1}, \dots, x_{j_k}$  respectively, are each of strictly smaller size than  $t$  as they are strict subterms of  $t$ .

In  $\ell_{\Gamma,y}[\ell_{\Delta,x}[\text{PI}(\pi)]]$ ,  $t$  is replaced by some  $z_j$ , which is existentially quantified. Hence  $t$  is a witness for  $z_j$  as due to the quantifier ordering, the existential quantification of  $z_j$  is in the scope of the quantifiers of  $x_{j_1}, \dots, x_{j_k}$  respectively. Therefore  $\Gamma \models Q_1 z_1 \dots Q_n z_n \ell_{\Gamma,y}[\ell_{\Delta,x}[\text{PI}(\pi)]]$ .

By Corollary 4.13  $\Delta \models \forall y_1 \dots \forall y_m \neg \ell_{\Gamma,y}[\text{PI}(\pi)]$ , where  $m$  is the number of  $\Gamma$ -colored terms in  $\text{PI}(\pi)$ . By a similar line of argumentation as above, we can replace the maximal  $\Delta$ -terms by existentially quantified variables and arrive at  $\Delta \models \overline{Q}_1 z_1 \dots \overline{Q}_n z_n \neg \ell_{\Delta,x}[\ell_{\Gamma,y}[\text{PI}(\pi)]]$  where  $\overline{Q}_i = \exists$  ( $\forall$ ) if  $Q_i = \forall$  ( $\exists$ ). Therefore also  $\Delta \models \neg Q_1 z_1 \dots Q_n z_n \ell_{\Delta,x}[\ell_{\Gamma,y}[\text{PI}(\pi)]]$ . By Lemma 4.14 and as all variables which replace colored terms are bound,  $\Delta \models \neg Q_1 z_1 \dots Q_n z_n \ell_{\Gamma,y}[\ell_{\Delta,x}[\text{PI}(\pi)]]$ .

As it is now easy to see that  $Q_1 z_1 \dots Q_n z_n \ell_{\Gamma,y}[\ell_{\Delta,x}[\text{PI}(\pi)]]$  contains no colored symbol, it is an interpolant.  $\square$

---

## Bibliography

- [BJ13] Maria Paola Bonacina and Moa Johansson. On Interpolation in Automated Theorem Proving. Technical Report 86/2012, Dipartimento di Informatica, Università degli Studi di Verona, 2013. Submitted to journal August 2013.
- [BL11] Matthias Baaz and Alexander Leitsch. *Methods of Cut-Elimination*. Trends in Logic. Springer, 2011.
- [CK90] C.C. Chang and H.J. Keisler. *Model Theory*. Studies in Logic and the Foundations of Mathematics. Elsevier Science, 1990.
- [Cra57a] William Craig. Linear Reasoning. A New Form of the Herbrand-Gentzen Theorem. *The Journal of Symbolic Logic*, 22(3):250–268, September 1957.
- [Cra57b] William Craig. Three Uses of the Herbrand-Gentzen Theorem in Relating Model Theory and Proof Theory. *The Journal of Symbolic Logic*, 22(3):269–285, September 1957.
- [DKPW10] Vijay D’Silva, Daniel Kroening, Mitra Purandare, and Georg Weissenbacher. Interpolant Strength. In *Proceedings of the International Conference on Verification, Model Checking, and Abstract Interpretation (VMCAI)*, volume 5944 of *Lecture Notes in Computer Science*, pages 129–145. Springer, January 2010.
- [Gen35] Gerhard Gentzen. Untersuchungen über das logische Schließen II. *Mathematische Zeitschrift*, 39, 1935.
- [Hua95] Guoxiang Huang. Constructing Craig Interpolation Formulas. In *Proceedings of the First Annual International Conference on Computing and Combinatorics, COCOON ’95*, pages 181–190, London, UK, UK, 1995. Springer-Verlag.
- [Kle67] Stephen Cole Kleene. *Mathematical logic*. Wiley, New York, NY, 1967.

- [Kra97] Jan Krajíček. Interpolation Theorems, Lower Bounds for Proof Systems, and Independence Results for Bounded Arithmetic. *Journal of Symbolic Logic*, pages 457–486, 1997.
- [Lyn59] Roger C. Lyndon. An Interpolation Theorem in the Predicate Calculus. *Pacific Journal of Mathematics*, 9(1):129–142, 1959.
- [McM03] Kenneth L. McMillan. Interpolation and SAT-Based Model Checking. In Jr. Hunt, Warren A. and Fabio Somenzi, editors, *Computer Aided Verification*, volume 2725 of *Lecture Notes in Computer Science*, pages 1–13. Springer Berlin Heidelberg, 2003.
- [Pud97] Pavel Pudlák. Lower Bounds for Resolution and Cutting Plane Proofs and Monotone Computations. *J. Symb. Log.*, 62(3):981–998, 1997.
- [Rob65] J. A. Robinson. A machine-oriented logic based on the resolution principle. *J. ACM*, 12(1):23–41, January 1965.
- [Sho67] Joseph R. Shoenfield. *Mathematical logic*. Addison-Wesley series in logic. Addison-Wesley Pub. Co., 1967.
- [Sla70] James R. Slagle. Interpolation theorems for resolution in lower predicate calculus. *J. ACM*, 17(3):535–542, July 1970.
- [Tak87] Gaisi Takeuti. *Proof Theory*. Studies in logic and the foundations of mathematics. North-Holland, 1987.
- [Wei10] Georg Weissenbacher. *Program Analysis with Interpolants*. PhD thesis, 2010.