

# Lecture Notes on Unification

15-317: Constructive Logic  
Frank Pfenning   André Platzer

Lecture 19  
April 9, 2020

In this lecture we take the essential step that makes the choice of goal and rule instantiation explicit in the operational semantics of logic programming. This consists of describing an algorithm for a problem called *unification* which, given two terms  $t$  and  $s$ , tries to find a substitution  $\theta$  for its free variables such that  $t\theta = s\theta$  if such a substitution exists. Recall that we write  $t\theta$  for the result of applying the substitution  $\theta$  to the term  $t$ .

These lecture notes are conceptual simplifications based on [1, 2].

## 1 Using Unification in Proof Search

Before we get to specifics of the algorithm, we consider how we use unification in proof search. Let us reconsider the (by now tired) example of unary addition

$$\frac{}{\text{plus}(0, N, N)} \text{pz} \qquad \frac{\text{plus}(M, N, P)}{\text{plus}(s(M), N, s(P))} \text{ps}$$

and an atomic goal such as

$$\text{plus}(s(0), s(s(0)), P).$$

Clearly the conclusion of the first rule does not match this goal, but the second one does. What question do we answer to arrive at this statement?

The first attempt might be: “*There is an instance of the rule such that the conclusion matches the goal.*” When we say *instance* we mean here the result of substituting terms for the variables occurring in a rule, proposition, or

term. We can see that this specification is not quite right: we need to instantiate the goal as well, since  $P$  must have the form  $s(P_1)$  for some as yet unknown  $P_1$ . The subgoal in that case would be  $\text{plus}(0, s(s(0)), P_1)$  according to the rule instantiation with 0 for  $M$ ,  $s(s(0))$  for  $N$ , and  $P_1$  for  $P$ .

The second attempt would therefore be: *“There is an instance of the rule and an instance of the goal such that the two are equal.”* This does not quite capture what we need either. For example, substituting  $s(s(s(s(0))))$  for  $P$  in the goal and  $s(s(s(0)))$  for  $P$  in the rule, together with the substitution for  $M$  and  $N$  from above, will also make the goal and conclusion of the rule identical, but is nonetheless wrong. The problem is that it would overcommit: using  $P_1$  for  $P$  in the rule, on the other hand, keeps the options open.  $P_1$  will be determined later by search and other unifications. In order to express this, we define that  $t_1$  is *more general* than  $t_2$  if  $t_1$  can be instantiated to  $t_2$ .

The third attempt is therefore: *“Find the most general instance of the rule and the goal so that the conclusion of the rule is equal to the instantiated goal.”* Phrased in terms of substitutions, this says: find  $\theta_1$  and  $\theta_2$  such that  $P'\theta_1 = P\theta_2$ , and any other common instance of  $P'$  and  $P$  is an instance of  $P'\theta_1$ .

In terms of the algorithm description it is more convenient if we redefine the problem slightly in this way: *“First rename the variables in the rule so that they are disjoint from the variables in the goal. Then find a single most general substitution  $\theta$  that unifies the renamed conclusion with the goal.”* Here, a unifying substitution  $\theta$  is most general if any other unifying substitution is an instance of  $\theta$ .

In the remainder of the lecture we will make these notions more precise and present an algorithm to compute a most general unifier. In the next lecture we show how to reformulate the operational semantics to explicitly use most general unifiers. This means that for the first time the semantics will admit free variables in goals.

**Notation:** In this lecture, we follow the general convention in logic that logical variables are called  $x, y, z$  while nullary function symbols are called  $a, b, c$ . Function symbols are called  $f, g, h$  while predicate symbols are called  $p, q, r$ . Recall that Prolog has the convention of starting variables with upper-case letters  $X, Y, Z$  while starting all function and predicate symbols including nullary ones with lower-case letters  $a, b, c, f, g, h, p, q, r$ .

## 2 Substitution

A substitution  $\sigma$  is a function that replaces variables by terms. When applied to a term,  $\sigma$  works in the canonical way, that is by uniformly replacing all the variables it affects in the term by their replacement but leaving the term unchanged otherwise. Different notations exist in the literature. In this lecture we adopt the notation  $t\sigma$  for the result of applying substitution  $\sigma$  to term  $t$ , instead of the more canonical  $\sigma(t)$ .

**Definition 1 (Substitution)** *A substitution is a homomorphism on terms and finitely supported, i.e., a function  $\sigma : \text{Term} \rightarrow \text{Term}$  on the set of terms such that*

$$\begin{aligned} f(t_1, \dots, t_n)\sigma &= f(t_1\sigma, \dots, t_n\sigma) && \text{for all function symbols } f \text{ and terms } t_i \\ \sigma &= id && \text{for almost all variables} \end{aligned}$$

*That is, the domain  $\text{dom}(\sigma) = \{x : x\sigma \neq x\}$  of all variables that are affected by  $\sigma$  is finite.*

In fact, the effect of a substitution is uniquely determined already by describing its effect on the finitely many variables that it affects. A substitution  $\sigma$  is, thus, represented by listing the replacements for each of the variables in  $\text{dom } \sigma$ . The most common notation for the substitution  $\sigma$  that, for example, replaces  $x$  by  $r$  and simultaneously replaces  $y$  by  $s$  and  $z$  by  $t$  and leaves all other variables unchanged since  $\text{dom}(\sigma) = \{x, y, z\}$ , is denoted:

$$(r/x, s/y, t/z)$$

This substitution is usually pronounced “ $r$  for  $x$ ,  $s$  for  $y$ , and  $t$  for  $z$ ”. It is enough to specify this substitution like that, because its effect on any other term is uniquely determined by Def. 1. Of course, it does not make sense to mention a replacement for  $x$  multiple times such as in  $(r/x, s/x)$  since that ill-formed representation does not give a function.

Substitutions can be extended to work also on logical formulas. For example,  $p(t)\sigma = p(t\sigma)$  for all predicate symbols  $p$ . But more attention is needed with quantifiers to avoid capture of variables such as when using the above substitution on the formula  $(\forall x p(x)) \wedge (\forall y q(y))$ . Quantifiers do not occur during proof search in Prolog, since that is all about unifying a goal  $p(t)$  with a corresponding clause  $p(s)$ , possibly with multiple arguments.

### 3 Composing Substitutions

Since substitutions are functions, they can be composed and it is easy to see that their composition is still a substitution.

**Theorem 1 (Composition of substitutions)** *The composition  $\sigma \circ \tau$  of substitution  $\sigma$  after substitution  $\tau$  is a substitution. For compatibility with the notation for applying substitutions, it is denoted  $\tau\sigma$ .*

$$t(\tau\sigma) = t(\sigma \circ \tau) = (t\tau)\sigma$$

**Proof:** Both  $\sigma$  and  $\tau$  only affect finitely many variables, and so will their composition.

$$x(\tau\sigma) = x(\sigma \circ \tau) = \sigma(\tau(x)) = (x\tau)\sigma$$

Since both  $\sigma$  and  $\tau$  leave function symbols unchanged and work homomorphically on the arguments, the same condition will be satisfied when applying  $\sigma$  after  $\tau$ .

$$\begin{aligned} f(t_1, \dots, t_n)(\tau\sigma) &= f(t_1(\tau\sigma), \dots, t_n(\tau\sigma)) \stackrel{\text{i.h.}}{=} f((t_1\tau)\sigma, \dots, (t_n\tau)\sigma) \\ &= f(t_1\tau, \dots, t_n\tau)\sigma = (f(t_1, \dots, t_n)\tau)\sigma \end{aligned}$$

□

A particularly useful type of substitutions are those that are *idempotent*, i.e.,  $\sigma \circ \sigma = \sigma$ , which means that applying the substitution twice has the same effect as applying it only once. Substitutions in which the variable they substitute occurs in its own replacement will not be idempotent so that things change again when they are applied multiple times. For example, applying  $(f(x)/x)$  twice will turn  $g(x)$  into  $g(f(f(x)))$  while applying it once will only give  $g(f(x))$ . Idempotent substitutions are easier to work with, because we do not have to pay attention how often we apply it to a term.

A representation of the composed substitution  $\tau\sigma$  can be computed from the representations of the substitutions  $\sigma$  and  $\tau$ . If  $\text{dom } \sigma \cap \text{dom } \tau = \emptyset$ , the composed representation is easier just by applying the subsequent substitution  $\sigma$  already to the replacement list of  $\tau$ :

$$\tau\sigma = (t\sigma/x : t/x \in \tau) \cup \sigma$$

where  $\cup$  denotes the union of the representations of a substitution. If  $\text{dom } \sigma$  and  $\text{dom } \tau$  overlap, then the replacements that  $\sigma$  does to variables of  $\text{dom } \tau$

will have to be dropped during  $\cup$ . Those variables  $x$  will already have been replaced by  $(\tau x)$  by the time  $\sigma$  is applied, so are already covered on the left. Fortunately, Prolog unification during proof search works on clauses that have already been renamed to use disjoint variable names for different clauses as well as different for each use of the same clause. Such disjointness conditions can be helpful. The set of all variables occurring in the substitution terms is the *codomain*  $\text{cod}(\theta) = \bigcup_{x \in \text{dom}(\theta)} \text{FV}(x\theta)$ . By  $\text{FV}(t)$  we denote the set of all free variables of term  $t$ , which are all variables in  $t$ .<sup>1</sup>

If  $\text{dom } \sigma \cap \text{dom } \tau = \emptyset$ , the above representation of the composition  $\tau\sigma$  can be obtained step by step when considering the elements of the list representing  $\tau$ , where  $(\cdot)$  is the empty substitution:

$$\begin{aligned} (t/x, \tau)\sigma &= (t\sigma/x, \tau\sigma) \\ (\cdot)\sigma &= \sigma \end{aligned}$$

## 4 Unifiers

**Definition 2 (Unifier)** *The substitution  $\sigma$  is a unifier for the terms  $s$  and  $t$  if  $s\sigma = t\sigma$ . Two terms  $s$  and  $t$  are called unifiable if a unifier  $\sigma$  exists.*

Unifiers are not unique. For example, both  $\sigma_1 = (h(a)/x, z/y)$  and  $\sigma_2 = (h(a)/x, y/z)$  unify

$$f(x, g(y)) \text{ and } f(h(a), g(z)) \quad (1)$$

Neither of those unifiers is inherently better than the other. In contrast,  $\sigma_3 = (h(a)/x, b/y, b/z)$  also unifies (1). But  $\sigma_3$  is worse than  $\sigma_1$  and  $\sigma_2$ , because  $\sigma_3$  is unnecessarily specific since it specializes both variable  $y$  and variable  $z$  to the nullary constant symbol  $b$ . After having applied  $\sigma_1$  to unify the terms in (1), we can still obtain the result of using the more special unifier  $\sigma_3$  by subsequently instantiating  $y$  and  $z$  to  $b$ .

**Definition 3 (Most-general unifier)** *A substitution  $\mu$  is a most-general unifier for the terms  $s$  and  $t$  iff  $\mu$  unifies  $s$  and  $t$  and*

*for all unifiers  $\sigma$  of  $s$  and  $t$  there is a substitution  $\sigma'$  such that  $\sigma = \mu\sigma' = \sigma' \circ \mu$*

*In other words, most-general unifiers are exactly the maximal elements with respect to the order:  $\sigma \prec \mu$  iff there is a  $\sigma'$  such that  $\sigma = \sigma' \circ \mu$ .*

---

<sup>1</sup>Representations of compositions  $\sigma\theta$  of representations  $\sigma$  and  $\theta$  of substitutions are conceptually particularly easy when  $\text{dom}(\theta) \cap \text{cod}(\theta) = \emptyset$  for all the relevant substitutions (which implies idempotence) and  $\text{dom}(\sigma) \cap (\text{dom}(\theta) \cup \text{cod}(\theta)) = \emptyset$ , since there is no overlapping effect in that case [1].

By this definition,  $\sigma_3$  is certainly not a most-general unifier for (1) but can still be represented as  $\sigma_3 = (b/y, b/z) \circ \sigma_1$ . Likewise,  $\sigma_2 = (y/z) \circ \sigma_1$ . Since that, indeed, holds for all other unifiers of (1), it turns out that  $\sigma_1$  is a most-general unifier for (1).

Most-general unifiers are still not unique, because  $\sigma_2$  also is a most-general unifier of (1) as well. After establishing that  $\sigma_1$  is a most-general unifier for (1), which requires more thought, it is easy to see that  $\sigma_2$  is a most-general unifier using  $\sigma_1 = (z/y) \circ \sigma_2$ . Every unifier  $\sigma$  already has a substitution  $\sigma'$  such that  $\sigma = \sigma' \circ \sigma_1$  since  $\sigma_1$  is a most-general unifier. Thus, every unifier  $\sigma$  also has a substitution  $\sigma' \circ (z/y)$  such that

$$(\sigma' \circ (z/y)) \circ \sigma_2 = \sigma' \circ ((z/y) \circ \sigma_2) = \sigma' \circ \sigma_1 = \sigma$$

Yet, when we have two most-general unifiers, they cannot be too far apart either. If we have two most-general unifiers  $\mu, \mu'$ , then we can get from  $\mu$  to  $\mu'$  by composing with some  $\sigma'$  ( $\mu' = \sigma' \circ \mu$ ), and we can, conversely, also get from  $\mu'$  to  $\mu$  with some  $\sigma$  ( $\mu = \sigma \circ \mu'$ ). Neither this  $\sigma$  nor this  $\sigma'$  can possibly substitute any variable by another term that isn't itself a variable. Suppose  $\sigma$  did replace some  $x$  by some function term  $f(e)$ , then  $f$  will always remain there (if  $x$  occurs in the input terms  $s$  and  $t$  at all) so the substitution  $\sigma'$  can never undo that effect of adding an  $f$ . It, thus, turns out that the most-general unifier is almost unique. The most-general unifier of a set of terms is unique up to variable renaming:

**Lemma 2** *If  $\mu, \mu'$  are most-general unifiers for the terms  $s$  and  $t$ , then there is a variable renaming  $\sigma$  such that  $\mu = \sigma \circ \mu'$ . A variable renaming is a substitution whose only effect is to replace variables by variables, not by arbitrary terms, and that, moreover, never renames two different variables to the same variable (which could never be undone by another substitution).*

## 5 Unification

As usual in this class, we present the algorithm to compute a most general unifier as a judgment, via a set of inference rules. The judgment has the form  $t \doteq s \mid \theta$ , where we think of  $t$  and  $s$  as inputs and a most general unifier  $\theta$  as the output. We read  $t \doteq s \mid \theta$  as “ $t$  unifies with  $s$  under  $\theta$ ”. In order to avoid the  $n$ -ary nature of the list of arguments, we will have an auxiliary judgment  $\mathbf{t} \doteq \mathbf{s} \mid \theta$  for sequences of terms  $\mathbf{t}$  and  $\mathbf{s}$ . Notions such as application of substitution are extended to sequences of terms in the obvious way. We use  $(\cdot)$  to stand for an empty sequence of terms (as

well as the empty substitution, which is a sequence of term and variable pairs). We will use boldface letters to stand for sequences of terms.

We first consider function terms and term sequences.

$$\frac{\mathbf{t} \doteq \mathbf{s} \mid \theta}{f(\mathbf{t}) \doteq f(\mathbf{s}) \mid \theta} \quad \frac{}{(\cdot) \doteq (\cdot) \mid (\cdot)} \quad \frac{t \doteq s \mid \theta_1 \quad \mathbf{t}\theta_1 \doteq \mathbf{s}\theta_1 \mid \theta_2}{(t, \mathbf{t}) \doteq (s, \mathbf{s}) \mid \theta_1\theta_2}$$

Observe how the most-general unifier  $\theta$  of  $t \doteq s$  from the left premise is already applied to the sequence on the right. Second, the cases for variables.

$$\frac{}{x \doteq x \mid (\cdot)} \quad \frac{x \notin \text{FV}(t)}{x \doteq t \mid (t/x)} \quad \frac{t = f(\mathbf{t}), x \notin \text{FV}(t)}{t \doteq x \mid (t/x)}$$

The condition that  $t = f(\mathbf{t})$  in the last rule (i.e., that the term  $t$  begins with a function symbol  $f$ ) ensures that it does not overlap with the rule for  $x \doteq t$ . The condition  $x \notin \text{FV}(t)$  in the last two rules also ensures that they do not overlap with the rule for  $x \doteq x$ . The condition that  $x \notin \text{FV}(t)$  is necessary because, for example, the two terms  $x$  and  $f(x)$  do not have unifier: no matter what, the substitution  $f(x)\theta$  will always have one more occurrence of  $f$  than  $x\theta$  and hence the two cannot be equal. Observe how this relates to the fact that the substitution  $f(x)/x$  that we would otherwise read of from the rule for  $x \doteq f(x)$  is not idempotent since its replacement of  $x$  mentions  $x$  again.

The other situations where unification fails is an equation of the form  $f(\mathbf{t}) = g(\mathbf{s})$  for  $f \neq g$ , and two sequences of terms of unequal length. The latter can happen if function symbols are overloaded at different arities, in which case failure of unification is the correct result.

Because of the so-called occurs check  $x \notin \text{FV}(t)$  whenever introducing a substitution  $t/x$ , the above rules will only ever create substitutions that are idempotent, which can be shown by a straightforward induction. Furthermore, the only place where substitutions are merged is the rule for  $(t, \mathbf{t}) \doteq (s, \mathbf{s}) \mid \theta_1\theta_2$ . By the occurs check, no variable in  $\text{dom}(\theta_1)$  will occur in its replacement. Consequently, no variable from  $\text{dom}(\theta_1)$  will occur in  $\mathbf{t}\theta_1$  or  $\mathbf{s}\theta_1$  so also not in the right premise  $\mathbf{t}\theta_1 \doteq \mathbf{s}\theta_1 \mid \theta_2$  nor its unifier  $\theta_2$ . Thus, the composition  $\theta_1\theta_2$  in the conclusion  $(t, \mathbf{t}) \doteq (s, \mathbf{s}) \mid \theta_1\theta_2$  will be for substitutions with disjoint domains  $\text{dom}(\theta_1) \cap \text{dom}(\theta_2) = \emptyset$  and is, thus, straightforward on their representations.

## 6 Soundness

There are a number of properties we would like to investigate regarding the unification algorithm proposed in the previous section. In today's lecture we will look at the first property of soundness, that is, we would like to show that the substitution  $\theta$  is indeed a unifier. If the judgment  $t \doteq s \mid \theta$  holds, then, indeed, applying the resulting substitution  $\theta$  to  $t$  gives the exact same term as applying the substitution  $\theta$  to  $s$  does.

**Theorem 3 (Soundness)** *If  $t \doteq s \mid \theta$  then  $t\theta = s\theta$ .*

**Proof:** Since the rule for  $f(\mathbf{t}) \doteq f(\mathbf{s}) \mid \theta$  reduces unification to unification of a sequence of terms, we need to generalize the induction hypothesis to cover the auxiliary unification judgment on term sequences.

(i) If  $t \doteq s \mid \theta$  then  $t\theta = s\theta$ .

(ii) If  $\mathbf{t} \doteq \mathbf{s} \mid \theta$  then  $\mathbf{t}\theta = \mathbf{s}\theta$ .

The proof proceeds by mutual induction on the structure of the deduction  $\mathcal{D}$  of  $t \doteq s$  and  $\mathcal{E}$  of  $\mathbf{t} \doteq \mathbf{s}$ . This means that if one judgment appears in the premiss of a rule for the other, we can apply the appropriate induction hypothesis.

In the proof below we will occasionally refer to *equality reasoning*, which refers to properties of equality in our mathematical language of discourse, not properties of the judgment  $t \doteq s$ . There are also some straightforward lemmas we do not bother to prove formally, such as  $t(s/x) = t$  if  $x \notin \text{FV}(t)$ .

**Case:**  $\mathcal{D} = \frac{\mathcal{E} \quad \mathbf{t} \doteq \mathbf{s} \mid \theta}{f(\mathbf{t}) \doteq f(\mathbf{s}) \mid \theta}$  where  $t = f(\mathbf{t})$  and  $s = f(\mathbf{s})$ .

$$\begin{aligned} \mathbf{t}\theta &= \mathbf{s}\theta \\ f(\mathbf{t})\theta &= f(\mathbf{s})\theta \end{aligned}$$

By i.h.(ii) on  $\mathcal{E}$   
By definition of substitution

**Case:**  $\mathcal{D} = \frac{}{(\cdot) \doteq (\cdot) \mid (\cdot)}$  where  $\mathbf{t} = \mathbf{s} = (\cdot)$  and  $\theta = (\cdot)$ .

$$(\cdot)\theta = (\cdot)\theta$$

By equality reasoning



**Case:**  $\mathcal{E} = \frac{\mathcal{D}_1 \quad \mathcal{E}_2}{(t_1, \mathbf{t}_2) \doteq (s_1, \mathbf{s}_2) \mid \theta_1 \theta_2}$  where  $\mathbf{t} = (t_1, \mathbf{t}_2)$  and  $\mathbf{s} = (s_1, \mathbf{s}_2)$   
and  $\theta = \theta_1 \theta_2$ .

$t_1 \theta_1 = s_1 \theta_1$	By i.h.(i) on $\mathcal{D}_1$
$(t_1 \theta_1) \theta_2 = (s_1 \theta_1) \theta_2$	By equality reasoning
$t_1(\theta_1 \theta_2) = s_1(\theta_2 \theta_2)$	By substitution composition (Theorem 1)
$(\mathbf{t}_2 \theta_1) \theta_2 = (\mathbf{s}_2 \theta_1) \theta_2$	By i.h.(ii) on $\mathcal{E}_2$
$\mathbf{t}_2(\theta_1 \theta_2) = \mathbf{s}_2(\theta_1 \theta_2)$	By substitution composition
$(t_1, \mathbf{t}_2)(\theta_1 \theta_2) = (s_1, \mathbf{s}_2)(\theta_1 \theta_2)$	By defn. of substitution

**Case:**  $\mathcal{D} = \frac{}{x \doteq x \mid (\cdot)}$  where  $t = s = x$  and  $\theta = (\cdot)$ .

$x(\cdot) = x(\cdot)$  By equality reasoning

**Case:**  $\mathcal{D} = \frac{x \notin \text{FV}(s)}{x \doteq s \mid (s/x)}$  where  $t = x$  and  $\theta = (s/x)$ .

$x(s/x) = s$  By defn. of substitution  
 $= s(s/x)$  Since  $x \notin \text{FV}(s)$

**Case:**  $\mathcal{D} = \frac{t = f(\mathbf{t}), x \notin \text{FV}(t)}{t \doteq x \mid (t/x)}$  where  $s = x$  and  $\theta = (t/x)$ .

$t(t/x) = t$  Since  $x \notin \text{FV}(t)$   
 $= x(t/x)$  By defn. of substitution

□

In a sense, soundness of the unification procedure is all that a correct deduction in logic programming depends on, since that is what determines whether a rule has been matched in plausible ways by a proper unifier. So whenever  $t \doteq s \mid \theta$ , then  $\theta$  is indeed a unifier of  $t$  and  $s$ . However, if  $\theta$  is a unifier but not a most-general unifier, then proof search may be incomplete, because  $\theta$  might then overcommit to an overly specific substitution for which no proof exists, while a proof would have existed for the most-general unifier (recall Section 1).

Prolog, however, actually prescribes a fixed proof search strategy: first applicable clause for the left-most subgoal is used with instantiation of their most-general unifier. If the unifier  $\theta$  from the judgment  $t \doteq s \mid \theta$  is merely a unifier but not most general, then we might not only have led to incomplete proof search, but might have also missed the fact that the subgoals unifies with another clause at all. Thus, for Prolog, we also need to make sure  $t \doteq s \mid \theta$  gives the most-general unifier and that we find out when  $t$  and  $s$  are not unifiable in the first place so that we can move on to consider the next clause instead. These questions will be considered in the next lecture.

## References

- [1] Frank Pfenning. Logic programming: Lecture 6: Unification. Lecture Notes 15-819K, Carnegie Mellon University, 2006.
- [2] Peter H. Schmitt. Formale Systeme. Vorlesungsskriptum Fakultät für Informatik , Universität Karlsruhe, 2012.