# Granularity-Adaptive Proof Presentation

Marvin Schiller
German Research Center for Artificial Intelligence
(DFKI), Bremen, Germany
Marvin.Schiller@dfki.de
and
Christoph Benzmüller
International University in Germany, Bruchsal, Germany
c.benzmueller@googlemail.com

# Granularity-Adaptive Proof Presentation

Marvin Schiller

German Research Center for Artificial Intelligence (DFKI), Bremen, Germany

`Marvin.Schiller@dfki.de`

and

Christoph Benzmüller

International University in Germany, Bruchsal, Germany

`c.benzmueller@googlemail.com`

February 27, 2009

**Abstract**

When mathematicians present proofs they usually adapt their explanations to their didactic goals and to the (assumed) knowledge of their addressees. Modern automated theorem provers, in contrast, present proofs usually at a fixed level of detail (also called granularity). Often these presentations are neither intended nor suitable for human use. A challenge therefore is to develop user- and goal-adaptive proof presentation techniques that obey common mathematical practice. We present a flexible and adaptive approach to proof presentation that exploits machine learning techniques to extract a model of the specific granularity of proof examples and employs this model for the automated generation of further proofs at an adapted level of granularity.

*Keywords*: Adaptive proof presentation, proof tutoring, automated reasoning, machine learning, granularity

## 1  Introduction

A key capability trained by students in mathematics and the formal sciences is the ability to conduct rigorous arguments and proofs and to present them. Thereby, proof presentation is usually highly adaptive as didactic goals and (assumed) knowledge of the addressee are taken into consideration. Modern automated theorem proving systems, however, do often not sufficiently address this common mathematical practice. They typically generate and present proofs using very fine-grained and machine-oriented calculi. While some theorem proving systems exist — amongst them prominent interactive theorem provers — that provide means for human-oriented proof presentations (e.g. proof presentation modules in Coq [17], Isabelle [18]) and Theorema [19]), the challenge of supporting user- and goal-adapted proof presentations has been widely neglected in the past. This constitutes an

[1] Let $x$ be an element of $A \cap (B \cup C)$, [2] then $x \in A$ and $x \in B \cup C$. [3] This means that $x \in A$, and either $x \in B$ or $x \in C$. [4] Hence we either have (i) $x \in A$ and $x \in B$, or we have (ii) $x \in A$ and $x \in C$. [5] Therefore, either $x \in A \cap B$ or $x \in A \cap C$, so [6] $x \in (A \cap B) \cup (A \cap C)$. [7] This shows that $A \cap (B \cup C)$ is a subset of $(A \cap B) \cup (A \cap C)$. [8] Conversely, let y be an element of $(A \cap B) \cup (A \cap C)$. [9] Then, either (iii) $y \in A \cap B$, or (iv) $y \in A \cap C$. [10] It follows that $y \in A$, and either $y \in B$ or $y \in C$. [11] Therefore, $y \in A$ and $y \in B \cup C$ so that $y \in A \cap (B \cup C)$. [12] Hence $(A \cap B) \cup (A \cap C)$ is a subset of $A \cap (B \cup C)$. [13] In view of Definition 1.1.1, we conclude that the sets $A \cap (B \cup C)$ and $(A \cap B) \cup (A \cap C)$ are equal.

Figure 1: Proof of the statement $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$, reproduced from [4]

unfortunate gap, in particular, since mathematics and the formal sciences are increasingly targeted as promising application areas for intelligent tutoring systems. In this paper we present a flexible and adaptive approach to proof presentation that exploits machine learning techniques to extract a model of the specific granularity of given proof examples, and that subsequently employs this model for the automated generation of further proofs at an adapted level of granularity. Our research has its roots in the collaborative DIALOG project [5] in which we developed means to employ the proof assistant ΩMEGA [16] for the dialog-based teaching of mathematical proofs. In DIALOG we have considered a dynamic approach: Instead of guiding the student along a pre-defined path towards a solution, we support the dynamic exploration of proofs, using automated proof search. This presupposes the development of techniques to adequately model the proofs a student is supposed to learn. Inference steps in ΩMEGA are implemented via an *assertion application* mechanism [8], which is based upon Serge Autexier's CORE calculus [1], as its logical kernel. In assertion level proofs, all inference steps are justified by a mathematical fact, such as definitions, theorems and lemmas, but not by steps of a purely technical nature such as structural decompositions, as required, for example, in natural deduction or sequent calculi.

The development of the dialog system prototype was guided by empirical studies using a mock-up of the DIALOG system [6]. One research challenge that educed out of the experiments is the question of judging the appropriate step size of proof steps (in the context of tutoring), also referred to as the *granularity* of mathematical proofs. Even in introductory textbooks in mathematics, intermediate proof steps are skipped, when this seems appropriate. An example is the elementary proof in basic set theory reproduced in Fig. 1. Whereas most of the proof steps consist of the application of exactly one mathematical fact (in this case, a definition or a lemma, such as the distributivity of *and* over *or*), the step from assertion [9] to assertion [10] suggests the application of several inference steps at once, namely the application of the definition of $\cap$ twice, and then using the distributivity of *and* over *or*.

Similar observations were made in the empirical studies within the DIALOG project. In these studies the tutors who helped to simulate the dialog system identified limits for how many inference steps are to be allowed at once. An example for an inacceptably large student step that was rejected by the tutor is presented right.

The idea to represent proofs at different levels of detail was incorporated into ΩMEGA

**student:** $(x, y) \in (R \circ S)^{-1}$
**tutor:**　　Now try to draw inferences from that!

| correct | appropriate | relevant |

**student:** $(x, y) \in S^{-1} \circ R^{-1}$
**tutor:**　　　　One cannot directly deduce that.

| correct | too coarse-grained | relevant |

as a hierarchically organized proof data structure [2]. The proof explanation system P.rex [9] implemented the idea to generate adapted proof presentations by moving up or down these layers on request. Alas, though the proofs at different level of detail can be handled by the ΩMEGA system, the problem remains of how to identify a particular level of granularity and how to ensure that this level of granularity is appropriate. This observation also applies to the Edinburgh HiProofs system [7].

Autexier and Fiedler have proposed one particular level of granularity [3], which they call *what-you-need-is-what-you-stated granularity*. Based on the assertion level inference mechanism in ΩMEGA, they also develop a proof checking mechanism for this level. In brief, their notion of granularity refers to such assertion level proofs, where all assertion level inference steps are spelled out explicitly and refer only to facts readily available from the assertions or the previous inference steps. However, they conclude that even the simple proof in Fig. 1 does not comply with their level of granularity, since the proof is missing some details.

This paper presents in Section 2 an adaptive framework to model proof granularity. This framework has been implemented as an extension of the ΩMEGA proof assistant and it is used to generate proof presentations at specific granularity levels of interest. In Section 3 we illustrate how our framework captures the granularity of our running example proof in Fig. 1. Models for granularity can be learned in our framework from samples, for which we employ standard machine learning techniques, as demonstrated in Section 4.

## 2 An Adaptive Model for Granularity

We treat the granularity problem as a classification task: given a proof step, representing one or several assertion applications, we judge it as either *appropriate*, *too big* or *too small*. As our feature space we employ several mathematical and logical aspects of proof steps, but also aspects of cognitive nature. For example, we keep track of the background knowledge of the user in a student model.

We illustrate our approach with an example proof step in Fig. 1: 10 is derived from 9 by applying the definition of ∩ twice, and then using the distributivity of *and* over *or*. In this step (which corresponds to multiple assertion level inference steps) we make the following observations: (i) involved are two concepts: def. of ∩ and distributivity of *and* over *or*, (ii) the total number of assertion applications is three, (iii) all involved concepts have been previously applied in the proof, (iv) all manipulations apply to a common part in 9, (v) the names of the applied concepts are not explicitly mentioned, and (vi) two of the assertion applications belong to *naive set theory* (def. of ∩) and one of them relates to the domain of propositional logic (distributivity).

These observations can be represented as a feature vector[1], where, in our example, the feature "distinct concepts" receives a value of "2", and so forth. We express our models for classifying granularity as rule sets, which associate specific combinations of feature values to a corresponding granularity verdict ("appropriate", "too big" or "too small"). These rule sets may be hand-authored by an expert or they may be learned from empirical data as

---

[1]Currently, we use around twenty features which are domain-independent, plus an indicator feature for each definition or lemma, and one indicator feature for each theory.

we show in Section 4. Our algorithm for granularity-adapted proof presentation takes two arguments, a granularity rule set and an assertion level proof[2] as generated by $\Omega$MEGA. Fig. 2 shows the assertion level proof generated by $\Omega$MEGA for our running example; this proof is represented as a tree (or acyclic graph) in sequent-style notation and the proof steps are ordered. Currently we only consider plain assertion level proofs, and do not assume any prior hierarchical structure or choices between proof alternatives (as possible in $\Omega$MEGA). Our algorithm performs an incremental categorization of steps in the proof tree (where $n = 0, \ldots, k$ denotes the ordered proof steps in the tree; initially $n$ is 1):

**while** there exists a proof step $n$
**do** evaluate the granularity of the compound proof step $n$ (i.e., the proof step consisting of all assertion level inferences performed after the last step labeled *appropriate* – or the beginning of the proof, if none exists yet) with the given rule set twice, (i) assuming that the involved concepts are mentioned in the presentation of the step (an "explanation"), and (ii) assuming only the resulting formula is displayed.

1. **if** $n$ is *appropriate* with explanation, **then** label $n$ as *appropriate*; set $n = n + 1$
2. **if** $n$ is *too small* with explanation, but *appropriate* without explanation **then** label $n$ as *appropriate without explanation*; set $n = n + 1$
3. **if** $n$ is *too small* both with and without explanation, **then** label $n$ as *too small*; set $n = n + 1$
4. **if** $n$ is *too big*, **then** label $n - 1$ as *appropriate* without explanation (i.e., consider the previous step as appropriate), unless $n - 1$ is labeled *appropriate* already or $n$ is the first step in the proof (in this special case label $n$ as *appropriate* with explanation, set $n = n + 1$)

**od**


We thereby obtain a proof tree with labeled steps (resp. labeled nodes) which differentiates between those nodes that are categorized as appropriate for presentation and those which are considered too fine-grained. Proof presentations are generated by walking through the tree[3], skipping the steps labeled *too small*[4].

When modeling granularity as a categorization problem, we have to test the hypothesis that the combination of features we devise is salient enough for the classification task. I.e., we have to determine whether steps within a class (i.e. "appropriate", "too big" and "too small") can indeed be fruitfully characterized by specific combinations of feature values, and distinguished from the feature values that characterize the two other classes. Our methodology for evaluation of this hypothesis consists in case studies and in empirical evaluations with mathematics tutors. This is exemplified in the following two sections.

---

[2]Our approach is not restricted to assertion level proofs and is also applicable to other calculi. However, in mathematics education we consider single assertion level proof steps as the finest granularity level of interest. We gained evidence for this choice from the empirical investigations in the DIALOG project (cf. [5] and [6]).

[3]In case of several branches, a choice is possible which subtree to present first, a question which we do not address in this paper.

[4]Even though the intermediate steps which are *too small* are withheld, the presentation of the output step reflects the results of all intermittent assertion applications, since we include the names of all involved concepts whenever a (compound) step is appropriate with explanation.

$$
\begin{array}{ll}
\textsc{Def}\cup\ (8) & \dfrac{x \in \mathbf{S} \vdash x \in \mathbf{S}}{} \\
\textsc{Def}\cap\ (7) & \dfrac{(x \in (A\cap B) \vee x \in (A\cap C)) \vdash x \in \mathbf{S}}{} \\
\textsc{Def}\cap\ (6) & \dfrac{(x \in (A\cap B) \vee x \in A \wedge x \in C) \vdash x \in \mathbf{S}}{} \\
\textsc{distr}(5) & \dfrac{(x \in A \wedge x \in B \vee x \in A \wedge x \in C) \vdash x \in \mathbf{S}}{} \\
\textsc{Def}\cup\ (4) & \dfrac{(x \in A \wedge (x \in B \vee x \in C)) \vdash x \in \mathbf{S}}{} \\
\textsc{Def}\cap\ (3) & \dfrac{(x \in A \wedge x \in (B\cup C)) \vdash x \in \mathbf{S}}{} \\
\textsc{Def}\subseteq\ (2) & \dfrac{(x \in (A \cap (B\cup C))) \vdash x \in \mathbf{S}}{} \\
\textsc{Def eq}\ (1) & \vdash (A\cap(B\cup C)) \subseteq \mathbf{S}
\end{array}
$$

$$
\begin{array}{ll}
\dfrac{y \in \mathbf{T} \vdash y \in \mathbf{T}}{} & \textsc{Def}\cap\ (15) \\
\dfrac{(y \in A \wedge y \in (B\cup C)) \vdash y \in \mathbf{T}}{} & \textsc{Def}\cup\ (14) \\
\dfrac{(y \in A \wedge (y \in B \vee y \in C)) \vdash y \in \mathbf{T}}{} & \textsc{distr}\ (13) \\
\dfrac{(y \in A \wedge y \in B \vee y \in A \wedge y \in C) \vdash y \in \mathbf{T}}{} & \textsc{Def}\cap\ (12) \\
\dfrac{(y \in A \wedge y \in B \vee y \in (A\cap C)) \vdash y \in \mathbf{T}}{} & \textsc{Def}\cap\ (11) \\
\dfrac{(y \in (A\cap B) \vee y \in (A\cap C)) \vdash y \in \mathbf{T}}{} & \textsc{Def}\cup\ (10) \\
\dfrac{(y \in \mathbf{S}) \vdash y \in \mathbf{T}}{} & \textsc{Def}\subseteq\ (9) \\
\vdash ((A \cap B) \cup (A \cap C)) \subseteq \mathbf{T}
\end{array}
$$

$$
\vdash \underbrace{(A\cap(B\cup C))}_{\mathbf{T}} = \underbrace{((A\cap B)\cup(A\cap C))}_{\mathbf{S}}
$$

Figure 2: Assertion level proof for the statement $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$

## 3 Case Study

In this section, we exemplarily model the step size of the textbook proof in Fig. 1. Starting point for the automated generation of our proof presentations are assertion level proofs in the mathematical assistant system $\Omega\textsc{mega}$. The basic assertion level proof, assuming the basic definitions in naive set theory, is presented in Fig. 2 as a sequent style proof tree. This proof consists of fifteen assertion level inference applications, which refer to the definitions of equality, subset, union and intersection as well as the concept of distributivity. Notice that the proof in Fig. 1 (taken from the textbook Bartle & Sherbert [4]) starts (in statement 1) with the assumption that an element $x$ is in the set $A \cap (B \cup C)$. The intention is to show the subset relation $A\cap(B\cup C) \subseteq (A\cap B)\cup(A\cap C)$. However, this is not explicitly revealed until step 6, when this part of the proof is already finished. The same style of *delayed* justification for prior steps is employed towards the end of the proof, where statements 12 and 13 justify (or recapitulate) the preceding proof. It must be questioned whether this style of presentation, where the motivation for some of the steps (such as the above assumption) is only presented in retrospective (when the assumption is discharged), is the still the most effective one for instructing students in our times. This style originated in former centuries, when the general task of the apprentice was to figure out the reason behind the procedures of his technically highly competent master with often poor teaching skills. Thus, for the modeling of step size, we consider a re-ordered variant of the steps in Fig. 1, which is displayed in Fig. 3 (a)[5]. We now generate a proof presentation which matches the step size of the twelve steps in the original proof, skipping intermediate proof steps according to our feature-based granularity model. Fig. 4 shows two sample rule sets which both lead to the proof presentation in Fig. 3 (b). The rule set in Fig. 4 (a) was generated by hand, whereas the rule set in Fig. 4 (b) was generated with the help of the C5 data mining tool [15][6]. The feature *hypintro* indicates whether a (multi-inference) proof step introduces a new

---

[5]Note that step (1) in the re-ordered proof corresponds to the statements 7, 12 and 13 in the original proof which jointly apply the concept of set equality.

[6]The sample proof was used to fit the rule set to it. All steps in the sample proof were provided as *appropriate*, all intermediate assertion level steps were labeled as *too-small*, and always the next bigger step to each step in the original proof was provided as an example for a *too big* step. Care was taken that the default rule of the generated rule set is of class *appropriate* (which was achieved via the cost function), so that the rule set better transfers to other domains. Otherwise, in case the default class was *too small*, and the examined proof steps were sufficiently different from the generating sample (and thus failed to match

1. In view of Definition 1.1.1, we [show] that the sets $A \cap (B \cup C)$ and $(A \cap B) \cup (A \cap C)$ are equal. ☐13 [First we show] that $A \cap (B \cup C)$ is a subset of $(A \cap B) \cup (A \cap C)$. ☐7 [Later we show] $(A \cap B) \cup (A \cap C)$ is a subset of $A \cap (B \cup C)$.☐12
2. Let $x$ be an element of $A \cap (B \cup C)$, ☐1
3. then $x \in A$ and $x \in B \cup C$. ☐2
4. This means that $x \in A$, and either $x \in B$ or $x \in C$.☐3
5. Hence we either have (i) $x \in A$ and $x \in B$, or we have (ii)$x \in A$ and $x \in C$.☐4
6. Therefore, either $x \in A \cap B$ or $x \in A \cap C$,☐5
7. so $x \in (A \cap B) \cup (A \cap C)$. ☐6
8. Conversely, let y be an element of $(A \cap B) \cup (A \cap C)$. ☐8
9. Then, either (iii) $y \in A \cap B$, or (iv) $y \in A \cap C$.☐9
10. It follows that $y \in A$, and either $y \in B$ or $y \in C$. ☐10
11. Therefore, $y \in A$ and $y \in B \cup C$, ☐11
12. so that $y \in A \cap (B \cup C)$. ☐11

(a)

1. We show that $((A \cap B) \cup (A \cap C) \subseteq A \cap B \cup C)$ and $(A \cap B \cup C \subseteq (A \cap B) \cup (A \cap C))$ ...because of definition of equality
2. We assume $x \in A \cap B \cup C$ and show $x \in (A \cap B) \cup (A \cap C)$
3. Therefore, $x \in A \wedge x \in B \cup C$
4. Therefore, $x \in A \wedge (x \in B \vee x \in C)$
5. Therefore, $x \in A \wedge x \in B \vee x \in A \wedge x \in C$
6. Therefore, $x \in A \cap B \vee x \in A \cap C$
7. We are done with the current part of the proof (i.e., to show that $x \in (A \cap B) \cup (A \cap C)$). [It remains to be shown that $(A \cap B) \cup (A \cap C) \subseteq A \cap B \cup C$]
8. We assume $y \in (A \cap B) \cup (A \cap C)$ and show $y \in A \cap B \cup C$
9. Therefore, $y \in A \cap B \vee y \in A \cap C$
10. Therefore, $y \in A \wedge (y \in B \vee y \in C)$
11. Therefore, $y \in A \wedge y \in B \cup C$
12. This finishes the proof. Q.e.d.

(b)

Figure 3: Comparison between (a) the (re-ordered) proof by Bartle and Sherbert [4] and (b) the proof presentation generated with our rule set from the $\Omega$MEGA proof in Fig. 2

hypothesis, and *close* indicates whether a branch of the proof has been finished. The feature *total* counts the number of assertion level inferences within one (multi-inference) step. Furthermore, the features *masteredconceptsunique* and *unmasteredconceptsunique* indicate how many of the employed concepts (if any) are supposed to be mastered resp. unmastered by the user according to a very basic student model (which is updated in the course of the proof). Furthermore, the occurrences of particular defined notions are counted (via the features ∩-Defn, ∪-Defn, equalitydefn). For example, the first rule in Fig. 4 (a) can be interpreted as 'if a step introduces a new hypothesis into the proof, and consists of more than one assertion level inference rule, it is considered too big'. Note that rules 4-6 in Fig. 4 (a) express the relation between the appropriateness of steps and whether the employed concepts are mentioned verbally (feature *verb*). Rule 6 has the effect of enforcing that the use of the definition of equality is always explicitly mentioned (as in step 1. in Fig 3 (b)). All other cases, which are not covered by the previous rules, are subject to a default rule. Rules are ordered by utility for conflict resolution.

The generated proof presentation in Fig. 3 (b) consists, similarly to the proof in Fig. 3 (a), of twelve steps.The three assertion level steps (11), (12) and (13) are combined into one single step from (9) to (10) in Fig. 3 (b). Natural language is produced via simple patterns (a more exciting natural language generation is possible with Fiedler's mechanisms [9], which is not the purpose of this paper).

The rule sets in Fig. 4 can be successfully reused for other examples in the domains as well. In Fig. 5, we present the resulting proof presentation when applying the rule set in Fig. 4 (a) to a different proof exercise (namely, a proof of the theorem $(A \cap B) \backslash C = A \cap (B \backslash C)$).

---

the non-default rules), the resulting proof presentation would be excessively short.

1) hypintro=1 $\wedge$ total> 1 $\Rightarrow$ step-too-big
2) $\cup$-Defn$\in\{1,2\}\wedge\cap$-Defn$\in\{1,2\}$ $\Rightarrow$ step-too-big
3) $\cap$-Defn< 3 $\wedge$ $\cup$-Defn=0 $\wedge$ masteredconceptsunique=1 $\wedge$ unmasteredconceptsunique=0 $\Rightarrow$ step-too-small
4) total<2 $\wedge$ verb=true $\Rightarrow$ step-too-small
5) masteredconceptsunique<3 $\wedge$ unmasteredconceptsunique=0 $\wedge$ verb=true $\Rightarrow$ step-too-small
6) equalitydefn>0 $\wedge$ verb=false $\Rightarrow$ step-too-big
7) _ $\Rightarrow$ step-appropriate

(a)

1) conceptsunique $\in\{0,1\}$ $\wedge$ equalitydefn=0 $\wedge$ verb=true $\Rightarrow$ step-too-small
2) hypintro=0 $\wedge$ equalitydefn=0 $\wedge$ $\cup$-Defn=0 $\wedge$ verb=true $\Rightarrow$ step-too-small
3) conceptsunique $\in\{2,3,4\}$ $\wedge$ $\cup$-Defn $\in\{1,2,3\}$ $\Rightarrow$ step-too-big
4) hypintro $\in\{1,2,3,4\}$ $\wedge$ conceptsunique $\in\{2,3,4\}$ $\Rightarrow$ step-too-big
5) unmasteredconceptsunique=0 $\wedge$ total $\in\{0,1,2\}$ $\cap$-Defn $\in\{1,2\}$ $\wedge$ close=false $\Rightarrow$ step-too-small
6) equalitydefn $\in\{1,2\}$ $\wedge$verb=false $\Rightarrow$ step-too-big
7) equalitydefn$\in\{1,2\}$ $\wedge$ verb=true $\Rightarrow$ step-appropriate
8) equalitydefn=0 $\wedge$ verb=false $\Rightarrow$ step-appropriate
9) _ $\Rightarrow$ step-appropriate

(b)

Figure 4: Rule sets employed in the running example: (a) rule set generated by hand, (b) rule set generated using C5 (ordered by the rules' confidence values)

1. We show that $((A \cap B)\backslash C \subseteq A \cap B\backslash C)$ and $(A \cap B\backslash C \subseteq (A \cap B)\backslash C)$ ...because of definition of equality
2. We assume $x \in A \cap B\backslash C$ and show $x \in (A \cap B)\backslash C$
3. Therefore, $x \in A \wedge x \in B\backslash C$
4. Therefore, $x \in A \wedge x \in B \wedge \neg(x \in C)$
5. We are done with the current part of the proof (i.e., to show that $x \in (A \cap B)\backslash C$). It remains to be shown that $(A \cap B\backslash C \subseteq A \cap B\backslash C$.
6. We assume $y \in (A \cap B)\backslash C$ and show $y \in A \cap B\backslash C$
7. Therefore, $y \in A \wedge y \in B \wedge \neg(y \in C)$ similarly to steps nr. (3 4)
8. This finishes the proof. Q.E.D. ... similarly to step nr. 7

Figure 5: Sample proof presentation generated via the rule set in Fig. 4 (a) for the theorem $(A \cap B)\backslash C = A \cap (B\backslash C)$

# 4  Learning from Empirical Data

Classification problems are a well-investigated topic in the machine learning community. There exist off-the-shelf tools that allow to learn classifiers (like our rule sets) from annotated examples (*supervised* learning). In our case, an expert annotates proof steps with the labels *appropriate*, *too small* or *too big*. Representing the proof steps in $\Omega$MEGA has the advantage that all the features of a particular proof step are computed in the background, and combined automatically with the expert's judgments as training instances for the learning algorithm. Currently, our algorithm calls the C5.0 data mining tools [15, 14] – which support the learning of decision trees and of rule sets – to obtain classifiers for granularity.

As part of an ongoing evaluation, we have conducted a study where a mathematician (with tutoring experience) judged the granularity of 135 proof steps. These steps were presented to him via an $\Omega$MEGA-assisted environment which computed the feature values for granularity classification in the background. The step size of proof steps presented to the expert was randomized, such that each presented step corresponded to one, two, or three assertion level inference steps. The presented proofs belonged to one exercise in naive

```
PART decision list
------------------

total <= 2 AND total > 0 AND parapos <= 0: appropriate (85.0/4.0)

total <= 2 AND unmasteredconceptsunique <= 0: step-too-small (11.0/2.0)

parapos <= 0 AND samesub <= 0: step-too-big (22.0/5.0)

unmasteredconceptsunique <= 1 AND hypintro <= 0: appropriate (9.0)

: step-too-big (8.0/2.0)
```

Figure 6: Empirically learnt rule set. The feature *parapos* indicates whether an inference has been applied only once in a proof situation where it could have been applied twice, in the same direction. The feature *samesub* indicates whether all inference applications within a (multi-inference) step apply to the same formula (and the same subparts thereof).

set theory and three different exercises about relations. We evaluated rule learning using C5.0 on our sample using 10 fold cross validation, which resulted in a mean percentage of correct classification of 84.6%, and $\kappa = 0.62$. We also used the PART classifier [10] included in the Weka suite[7], which is inspired by Quinlan's C4.5. After we excluded some of the attributes (in particular those that refer to the use of specific concepts, i.e., Def. of $\cap$, Def. of $\circ$, etc...), PART achieved 86.7% of correctly classified instances in stratified cross validation ($\kappa$=0.68). Apparently, removal of the most domain-specific attributes prevented the algorithm from overfitting. The resulting rule set is presented in Fig. 6.

The feature *parapos* indicates whether an inference has been applied only once in a proof situation where it could have been applied twice, in the same direction. The feature *samesub* indicates whether all inference applications within a (multi-inference) step apply to the same formula (and the same subparts thereof). When applied to our running example, we obtain the proof presentation as shown in Fig. 7.

To compare the rule-based classifiers with support vector machines, we applied SMO [13] on our data, resulting in 83.0% correctness and $\kappa$=0.57 in stratified cross validation, which is a similar performance to C5.0.

# 5   Conclusion

Granularity has been a challenge in AI for decades [11, 12]. Here we have focused on adaptive proof granularity, which we treat as a classification problem. We model different levels of granularity using rule sets, which can be hand coded or learned from sample proofs.

As a case study, we have formulated the granularity level of the proof in Fig. 1 from the textbook [4] as a rule set in our classification-based approach. Classifiers are applied dynamically to each proof step, thus taking into account changeable information such as the user's familiarity with the involved concepts. Using assertion level proofs as the basis for our

---

[7]http://www.cs.waikato.ac.nz/~ml/weka/

1. We show that $((A \cap B) \vee (A \cap C) \subseteq A \cap B \vee C)$ and $(A \cap B \vee C \subseteq (A \cap B) \vee (A \cap C))$ ...because of definition of equality
2. We assume $x \in A \cap B \vee C$ and show $x \in (A \cap B) \vee (A \cap C)$ ...because of definition of subset
3. Therefore, $x \in A \wedge x \in B \vee C$ ...because of definition of intersection
4. Therefore, $x \in A \wedge (x \in B \vee x \in C)$ ...because of definition of union
5. Therefore, $x \in A \wedge x \in B \vee x \in A \wedge x \in C$ ...because of logics
6. Therefore, $x \in A \cap B \vee x \in A \cap C$ ...because of definition of intersection ... similarly to step nr. 3
7. We are done with the current part of the proof (i.e., to show that $x \in (A \cap B) \vee (A \cap C)$). It remains to be shown that $(A \cap B) \vee (A \cap C) \subseteq A \cap B \vee C$. ... because of definition of union.
8. We assume $y \in (A \cap B) \vee (A \cap C)$ and show $y \in A \cap B \vee C$ ...because of definition of subset
9. Therefore, $y \in A \cap B \vee y \in A \cap C$ ...because of definition of union
10. Therefore, $y \in A \wedge y \in B \vee y \in A \wedge y \in C$ ...because of definition of intersection ... similarly to step nr. 3
11. Therefore, $y \in A \wedge (y \in B \vee y \in C)$ ...because of logics
12. Therefore, $y \in A \wedge y \in B \vee C$ ...because of definition of union
13. This finishes the proof. Q.e.d. ...because of definition of intersection

Figure 7: The assertion level proof in Fig. 2 presented according to the rule set from Fig. 6

approach has the additional advantage that the relevant information for the classification task (e.g. the concept names) is easily read off the proofs. This also eases the generation of natural language proof output in general.

Future work consists in empirical evaluations of the learning approach, in order to address the following questions;

(i) what are the most salient features for judging granularity, and are they different among distinct experts and domains,

(ii) what is the interrater reliability among different experts and the corresponding classifiers generated by learning in our framework?

The resulting corpora of annotated proof steps and generated classifiers can then be used to evaluate the appropriateness of the proof presentations generated by our system.

# References

[1] S. Autexier. The CoRe calculus. Proc. *CADE-20*, vol. 3632 of *LNCS*, Springer, 2005.

[2] S. Autexier, C. Benzmüller, D. Dietrich, A. Meier, and C.-P. Wirth. A generic modular data structure for proof attempts alternating on ideas and granularity. Proc. *MKM'05*, vol. 3863 of *LNCS*, Springer, 2006.

[3] S. Autexier and A. Fiedler. Textbook proofs meet formal logic - the problem of underspecification and granularity. Proc. *MKM-05*, vol. 3863 of *LNCS*, Springer, 2005.

[4] R. G. Bartle and D. Sherbert. *Introduction to Real Analysis*. Wiley, 2nd edition, 1982.

[5] C. Benzmüller, H. Horacek, I. Kruijff-Korbayová, M. Pinkal, J. H. Siekmann, and M. Wolska. Natural language dialog with a tutor system for mathematical proofs. In *Cognitive Systems*, vol. 4429 of *LNAI*, Springer, 2005.

[6] C. Benzmüller, H. Horacek, H. Lesourd, I. Kruijff-Korbajova, M. Schiller, and M. Wolska. A corpus of tutorial dialogs on theorem proving; the influence of the presentation of the study-material. In *Proc. of Int. Conf. on Language Resources and Evaluation (LREC 2006)*, Genoa, Italy, 2006. ELDA.

[7] E. Denney, J. Power, and K. Tourlas. Hiproofs: A hierarchical notion of proof tree. Electr. Notes Theor. Comput. Sci., 155: 341-359, 2006.

[8] D. Dietrich. The task-layer of the ΩMEGA system. Diploma thesis, FR 6.2 Informatik, Universität des Saarlandes, Saarbrücken, Germany, 2006.

[9] A. Fiedler. *P.rex*: An interactive proof explainer. Proc. *IJCAR 2001*, vol. 2083 of LNAI, Springer, 2001.

[10] E. Frank and I. H. Witten. Generating accurate rule sets without global optimization. In *Proc. of 15th Int. Conf. on Machine Learning*. Morgan Kaufmann, 1998.

[11] J. R. Hobbs. Granularity. Proc. *IJCAI-9*, pp. 432–435, Los Angeles, CA, USA, 1985.

[12] G. Mccalla, J. Greer, B. Barrie, and P. Pospisil. Granularity hierarchies. In *Computers and Mathematics with Applications: Sp. Issue on Semantic Networks*, pp. 363–375, 1992.

[13] J. Platt. Fast Training of Support Vector Machines Using Sequential Minimal Optimization. In B. Schoelkopf, C. Burges, and A. Smola, editors, *Advances in Kernel Methods - Support Vector Learning*, pp. 185-208. MIT Press, 1999.

[14] J. R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, 1993.

[15] RuleQuest Research. Data mining tools see5 and c5.0. `http://www.rulequest.com/see5-info.html`, 2007.

[16] J. H. Siekmann, C. Benzmüller, and S. Autexier. Computer supported mathematics with omega. *Journal of Applied Logic*, 4(4):533–559, 2006.

[17] Y. Coscoy, G. Kahn, L. Thery. Extracting text from proofs. In *Proc. Typed Lambda Calculus and its Applications*, pp. 109-123, Edinburgh, UK, Springer, 1995.

[18] M. Simons. Proof Presentation for Isabelle. In *Proc. TPHOLs '97*, London, UK, Springer, 1997.

[19] B. Buchberger, T. Jebelean, F. Kriftner, M. Marin, E. Tomuta, and D. Vasaru. A Survey of the Theorema project. In *Proc. ISSAC'97*, Maui, Hawaii, pp. 384–391, 1997.