

CASE-16

CASE-16

CASE-16

CASE-16

Proceedings of the 6th IJCAR ATP System Competition (CASC-J6)

Geoff Sutcliffe

University of Miami, USA

Abstract

The CADE ATP System Competition (CASC) evaluates the performance of sound, fully automatic, classical logic, ATP systems. The evaluation is in terms of the number of problems solved, the number of acceptable proofs and models produced, and the average runtime for problems solved, in the context of a bounded number of eligible problems chosen from the TPTP problem library, and specified time limits on solution attempts. The 6th IJCAR ATP System Competition (CASC-J6) was held on 24th and 28th June 2012. The design of the competition and its rules, and information regarding the competing systems, are provided in this report.

1 Introduction

The CADE conferences are the major forum for the presentation of new research in all aspects of automated deduction. In order to stimulate ATP research and system development, and to expose ATP systems within and beyond the ATP community, the CADE ATP System Competition (CASC) is held at each CADE conference. CASC-J6 was held on 24th and 28th June 2012, as part of the Alan Turing Centenary Conference and the 6th International Joint Conference on Automated Reasoning (IJCAR 2012)¹, in Manchester, England. It is the seventeenth competition in the CASC series [117, 122, 120, 88, 90, 116, 114, 115, 95, 97, 99, 101, 104, 107, 109, 110].

CASC evaluates the performance of sound, fully automatic, classical logic, ATP systems. The evaluation is in terms of:

- the number of problems solved,
- the number of acceptable proofs and models produced, and
- the average runtime for problems solved;

in the context of:

- a bounded number of eligible problems, chosen from the TPTP problem library [105], and
- specified time limits on solution attempts.

Thirty-five ATP systems and variants, listed in Tables 1 and 2 entered into the various competition and demonstration divisions. The winners of the CASC-23 (the previous CASC) divisions were automatically entered into those divisions, to provide benchmarks against which progress can be judged (the competition archive provides access to the systems' executables and source code).

The design and procedures of this CASC evolved from those of previous CASCs [117, 118, 113, 119, 86, 87, 89, 91, 92, 93, 94, 96, 98, 100, 103, 106, 108]. Important changes for this CASC were:

¹CADE was a constituent conference of IJCAR, hence CASC-“J6”.

ATP System	Divisions	Entrants (Associates)	Entrant's Affiliation
CVC4 0.0	FNN@T FNT EPR	Andrew Reynolds (Cesare Tinelli, Clark Barrett) Björn Pelzer	University of Iowa
E-Darwin 1.5	FOF@T FNT@T FOF FNT EPR	Björn Pelzer	University Koblenz-Landau
E-KRHyper 1.3	FOF@T FNT@T MZR@T FOF FNT EPR LTB FOF@T FOF LTB	Björn Pelzer	University Koblenz-Landau
E-Males 1.1	FOF	Daniel Kuehlwein (Josef Urban, Stephan Schulz) CASC	Radboud University Nijmegen
EP 1.4pre	FOF@T FNT@T MZR@T	Stephan Schulz	CASC-23 CNF winner
EP(-SAT) 1.6pre	FOF FNT EPR LTB		Technische Universität München
FIMO 0.3	FNT@T FNT	Orkunt Sabuncu	University of Potsdam
iProver 0.9	EPR	CASC	CASC-23 EPR winner
iProver(-SAT) 0.99	FOF@T FNT@T MZR@T	Konstantin Korovin (Christoph Stickel)	University of Manchester
iProver-Eq 0.8	FOF FNT EPR LTB FOF@T FNT@T MZR@T FOF FNT EPR LTB	Christoph Stickel (Konstantin Korovin)	University of Manchester
Isabelle 2012	THF	Jasmin Blanchette (Lawrence Paulson, Tobias Nipkow, Makarius Wenzel)	Technische Universität München
Isabelle+HO 2012	THF (demo)	Jasmin Blanchette (Lawrence Paulson, Tobias Nipkow, Makarius Wenzel)	Technische Universität München
leanCoP 2.2	FOF@T FOF	Jens Otten	University of Potsdam
leanCoP-ARDE 2.2	MZR@T LTB	Mario Frank (Thomas Rath, Jens Otten)	University of Potsdam

Table 1: The ATP systems and entrants

ATP System	Divisions	Entrants (Associates)	Entrant's Affiliation
LEO-II 1.4	THF FOF	Christoph Benz Müller	Free University Berlin
MaLAREa 0.4	MZR@T	Josef Urban (Daniel Kuehlwein, Stephan Schulz, Jiri Vyskocil)	Radboud University Nijmegen
Muscadet 4.2	FOF@T FOF	Dominique Pastre	University Paris Descartes
Nitrox 2012	FNT@T FNT	Jasmin Blanchette (Emina Torlak)	Technische Universität München
Paradox 3.0	FNT	CASC	CASC-23 FNT winner
Princess 2012-05-28	FOF@T TFA FOF	Philipp Ruegger (Aleksandar Zeljic)	Uppsala University
Prover9 2009-11A	FOF	CASC (William McCune, Bob Veroff)	CASC fixed point
PS-E 1.0	MZR@T	Daniel Kuehlwein (Josef Urban, Stephan Schulz)	Radboud University Nijmegen
Satallax 2.1	THF	CASC	CASC-23 THF winner
Satallax 2.4	THF	Chad Brown	Saarland University
SPASS+T 2.2.14	TFA	CASC	CASC-23 TFA winner
SPASS+T 2.2.16	TFA	Uwe Waldmann (Stephan Zimmer)	Max-Planck-Institut für Informatik
STP 1.0	FOF@T FNT@T FOF FNT EPR	Adam Pease (Stephan Schulz)	Articulate Software
SuperZenon 0.0.1	FOF@T FOF	David Delahaye (Melanie Jacquiel)	CEDRIC/CNAM
TPS 3.120601S1b	THF	Chad E. Brown (Peter Andrews)	Saarland University
Vampire 0.6	FOF	CASC	CASC-23 FOF winner
Vampire-LTB 1.8	LTB	CASC	CASC-23 LTB winner
Vampire 2.6	FOF@T FNT@T MZR@T	Krystof Hoder (Andrei Voronkov)	University of Manchester
Zenon 0.7.1	FOF FNT EPR LTB FOF@T FOF	Damien Doligez	INRIA

Table 2: The ATP systems and entrants, continued

- CASC-J6 was run in two parts: CASC@Turing was held during the Alan Turing Centenary Conference, and the “regular” CASC was held during IJCAR. CASC@Turing had three divisions: FOF@Turing, FNT@Turing, MZR@Turing. The regular CASC had six divisions: THF, TFA, FOF, FNT, EPR, LTB.
- The FOF division included the previous CNF division. The FOF division had four problem categories: FEQ, FNE, CEQ, CNE.
- The FNT division included the previous SAT division (which had been suspended since CASC-J5). The FNT division had four problem categories: FNQ, FNN, SEQ, SNE.
- The UEQ division was suspended.
- Prover9 was entered as a replacement for Otter, as a fixed point against which progress can be judged over the years.

The competition organizer was Geoff Sutcliffe. The competition was overseen by a panel of knowledgeable researchers who were not participating in the event; the CASC-23 panel members were Franz Baader, Dale Miller, and Christoph Weidenbach. The CASC rules, specifications, and deadlines are absolute. Only the panel has the right to make exceptions. The competition was run on computers provided by the Max-Planck-Institut für Informatik, Saarbrücken, Germany. The CASC-J6 web site provides access to resources used before, during, and after the event: <http://www.tptp.org/CASC/J6>

It is assumed that all entrants have read the web pages related to the competition, and have complied with the competition rules. Non-compliance with the rules could lead to disqualification. A “catch-all” rule is used to deal with any unforeseen circumstances: *No cheating is allowed*. The panel is allowed to disqualify entrants due to unfairness, and to adjust the competition rules in case of misuse.

2 Divisions

CASC is run in divisions according to problem and system characteristics. There are *competition* divisions in which systems are explicitly ranked, and a *demonstration* division in which systems demonstrate their abilities without being formally ranked. Some divisions are further divided into problem categories, which make it possible to analyse, at a more fine grained level, which systems work well for what types of problems. The problem categories have no effect on the competition rankings, which are made at only the division level.

2.1 The Competition Divisions

The competition divisions are open to ATP systems that meet the required system properties, described in Section 6.1. Each competition division uses problems that have certain logical, language, and syntactic characteristics, so that the ATP systems that compete in the division are, in principle, able to attempt all the problems in the division.

The **FOF@Turing** division: First-Order Form syntactically non-propositional theorems (axioms with a provable conjecture). The FOF@Turing division has three problem categories:

- The **FNE** category: FOF with No Equality, not (known to be) effectively propositional
- The **FEQ** category: FOF with Equality, not (known to be) effectively propositional
- The **FEP** category: FOF Effectively Propositional

The **FNT@Turing** division: First-order form syntactically non-propositional Non-Theorems (axioms with an unprovable conjecture, and satisfiable axioms sets). The FNT@Turing division has two problem categories:

- The **FNN** category: FNT with No equality
- The **FNQ** category: FNT with eQuality

The **MZR@Turing** division: First-order form non-propositional theorems (axioms with a provable conjecture) taken from an export of Mizar problems, presented in batches. The problems in each batch can be attempted in any order, including multiple attempts on a problem. This facilitates learning from proofs found. There is no per-problem time limit, only an overall wall clock time limit. A set of training problems and solutions was available for tuning and training systems entered in this division.

The **THF** division: Typed Higher-order Form non-propositional theorems (axioms with a provable conjecture), using the THF0 syntax. The THF division has two problem categories:

- The **TNE** category: THF with No Equality
- The **TEQ** category: THF with EQuality

The **TFA** division: Typed First-order with Arithmetic theorems (axioms with a provable conjecture, using the TFF0 syntax. The TFA division has two problem categories:

- The **TFI** category: TFA with only Integer arithmetic
- The **TFR** category: TFA with only Rational and only Real arithmetic (no mixed rational and real arithmetic)

The **FOF** division: First-Order Form syntactically non-propositional theorems (axioms with a provable conjecture), and clause normal form really non-propositional theorems (unsatisfiable clause sets), but not unit equality problems. The FOF division has four problem categories:

- The **FNE** category: FOF with No Equality effectively propositional
- The **FEQ** category: FOF with EQuality effectively propositional
- The **CNE** category: CNF with No Equality
- The **CEQ** category: CNF with EQuality

The **FNT** division: First-order form syntactically non-propositional Non-Theorems (axioms with a countersatisfiable conjecture, and satisfiable axiom sets), and clause normal form really non-propositional non-theorems (satisfiable clause sets). The FNT division has four problem categories:

- The **FNN** category: FNT with No equality
- The **FNQ** category: FNT with eQuality
- The **SNE** category: SAT with No Equality
- The **SEQ** category: SAT with EQuality

The **EPR** division: Effectively PRopositional clause normal form theorems and non-theorems (clause sets). *Effectively propositional* means non-propositional with a finite Herbrand Universe. The EPR division has two problem categories:

- The **EPT** category: Effectively Propositional Theorems (unsatisfiable clause sets)
- The **EPS** category: Effectively Propositional non-theorems (Satisfiable clause sets)

The **LTB** division: First-order form non-propositional theorems (axioms with a provable conjecture) from Large Theories, presented in Batches. The LTB division has three problem categories:

- The **ISA** category: Problems exported from Isabelle.
- The **MZR** category: Problems exported from the Mizar Mathematical Library (MML).
- The **SMO** category: Problems exported from the Suggested Upper Merged Ontology (SUMO).

Section 3.2 explains what problems are eligible for use in each division and category. Section 4 explains how the systems are ranked in each division.

2.2 The Demonstration Division

ATP systems that cannot run in the competition divisions for any reason (e.g., the system requires special hardware, or the entrant is an organizer) can be entered into the demonstration division. Demonstration division systems can run on the competition computers, or the computers can be supplied by the entrant. Computers supplied by the entrant may be brought to CASC, or may be accessed via the internet. The demonstration division results are presented along with the competition divisions' results, but might not be comparable with those results. The systems are not ranked and no prizes are awarded.

3 Infrastructure

3.1 Computers

The computers were Dell PowerEdge blade computers, each having:

- Two Intel Xeon E5620, quad-core, 2.40GHz CPUs, hyperthreading enabled
- 48GB memory
- GNU Linux c15-001 2.6.32.22.1.amd64-smp operating system

In the non-batch divisions systems could use only one core, and were limited to a fraction of the memory, as explained in Section 3.3 (multiple jobs were run on each node). In the batch divisions each system was allocated one node, and could use all the cores and memory.

3.2 Problems

3.2.1 Problem Selection

The problems were taken from the TPTP problem library, version v5.4.0. Additionally, problems for the MZR@Turing division were taken from the MPTP2078 problem set², and problems for the SMO category of the LTB division were taken from a problem set developed for CASC-J6. The TPTP version used for CASC is released after the competition has started, so that new problems have not been seen by the entrants. Access to and use of the non-TPTP problem sets was controlled to ensure that the system complied with the CASC tuning restrictions, described in Section 6.1.

The problems have to meet certain criteria to be eligible for selection:

²<http://wiki.mizar.org/twiki/bin/view/Mizar/MpTP2078>

- The TPTP uses system performance data to compute problem difficulty ratings [121], and from the ratings classifies problems as one of:
 - Easy: Solvable by all state-of-the-art ATP systems
 - Difficult: Solvable by some state-of-the-art ATP systems
 - Unsolved: Solvable by no ATP systems
 - Open: Theoremhood unknown

Difficult problems with a rating in the range 0.21 to 0.99 are eligible. Problems of lesser and greater ratings might also be eligible in some divisions (especially the batch divisions, because the TPTP problem ratings are computed from non-batch mode results). Performance data from systems submitted by the system submission deadline is used for computing the problem ratings for the TPTP version used for the competition.

- The TPTP distinguishes versions of problems as one of standard, incomplete, augmented, especial, or biased. All except biased problems are eligible.
- In the batch divisions, the problems are selected so that there is consistent symbol usage between problems in each category, but there may not be consistent axiom naming between problems.

The problems used are randomly selected from the eligible problems at the start of the competition, based on a seed supplied by the competition panel.

- The selection is constrained so that no division or category contains an excessive number of very similar problems.
- The selection mechanism is biased to select problems that are new in the TPTP version used, until 50% of the problems in each category have been selected, after which random selection (from old and new problems) continues. The actual percentage of new problems used depends on how many new problems are eligible and the limitation on very similar problems.

3.2.2 Number of Problems

The minimal numbers of problems that must be used in each division and category, to ensure sufficient confidence in the competition results, are determined from the numbers of eligible problems in each division and category [44] (the competition organizers have to ensure that there are sufficient computers available to run the ATP systems on this minimal number of problems). The minimal numbers of problems is used in determining the time limits imposed on each solution attempt - see Section 3.3.

A lower bound on the total number of problems to be used is determined from the number of computers available, the time allocated to the competition, the number of ATP systems to be run on the competition computers over all the divisions, and the per-problem time limit, according to the following relationship:

$$NumberOfProblems = \frac{NumberOfComputers * TimeAllocated}{NumberOfATPSystems * TimeLimit}$$

It is a lower bound on the total number of problems because it assumes that every system uses all of the time limit for each problem. Since some solution attempts succeed before the time limit is reached, more problems can be used.

The numbers of problems used in the categories in the various divisions are (roughly) proportional to the numbers of eligible problems than can be used in the categories, after taking into account the limitation on very similar problems. The numbers of problems used in each division and category are determined according to the judgement of the competition organizers.

3.2.3 Problem Preparation

The problems are in TPTP format, with `include` directives (included files are found relative to the TPTP environment variable). The problems in each non-batch division, and each LTB batch, are given in increasing order of TPTP difficulty rating. This is aesthetic in the non-batch divisions, but practically important in the batches where it is possible to learn from proofs found earlier in the batch.

In order to ensure that no system receives an advantage or disadvantage due to the specific presentation of the problems in the TPTP, the problems are preprocessed to:

- strip out all comment lines, including the problem header
- randomly reorder the formulae/clauses (the `include` directives are left before the formulae, type declarations and definitions are kept before the symbols' uses)
- randomly swap the arguments of associative connectives, and randomly reverse implications
- randomly reverse equalities

In order to prevent systems from recognizing problems from their file names, symbolic links are made to the selected problems, using names of the form `CCCNNN-1.p` for the symbolic links. CCC is the division or problem category name, and NNN runs from 001 to the number of problems in the respective division or problem category. The problems are specified to the ATP systems using the symbolic link names.

In the demonstration division the same problems are used as for the competition divisions, with the same preprocessing applied. However, the original file names can be retained for systems running on computers provided by the entrant.

3.2.4 Batch Specification Files

The problems for each batch division and category are listed in a batch specification file, containing one or more *batch specifications*. Each batch specification consists of:

- A header line `% SZS start BatchConfiguration`
- A problem category line of the form
`division.category division_mnemonic.category_mnemonic`
 For the MZR@Turing division it was
`division.category MZR.MZR`
 For the LTB division it was
`division.category LTB.category_mnemonic`
- A specification of whether or not the problems in the batch must be attempted in order is given, in a line of the form
`execution.order ordered/unordered`
 For the MZR@Turing division it was
`execution.order unordered`
 For the LTB division it was
`execution.order ordered`
 i.e., in the LTB division systems may not start any attempt on a problem, including reading the problem file, before ending the attempt on the preceding problem.
- A specification of what output is required from the ATP systems for each problem, in a line of the form
`output.required space_separated_list`

where the available list values are the SZS values **Assurance**, **Proof**, **Model**, and **Answer**. For CASC-J6 it was

```
output.required Assurance.
```

- A specification of what output is desired from the ATP systems for each problem, in a line of the form

```
output.desired space_separated_list
```

where the list values are as for the required output. For the MZR@Turing division and the MZR problem category of the LTB division it was

```
output.desired Proof
```

For the ISA problem category it was

```
output.desired Proof ListOfFOF
```

where the **ListOfFOF** is a list of axioms sufficient for a proof (a well-formed proof can be used to provide such a list). For the SMO problem category it was

```
output.desired Proof Answer
```

where the answer is a definite binding for the outermost existentially quantified variables of the conjecture.

- The wall clock time limit per problem, in a line of the form

```
limit.time.problem.wc limit_in_seconds
```

A value of zero indicates no per-problem limit - this will be the case for the MZR@Turing division.

- A terminator line **% SZS end BatchConfiguration**
- A header line **% SZS start BatchIncludes**
- **include** directives that are used in every problem. Problems in the batch have all these **include** directives, and can also have other **include** directives that are not listed here.
- A terminator line **% SZS end BatchIncludes**
- A header line **% SZS start BatchProblems**
- Pairs of absolute problem file names, and absolute output file names where the output for the problem must be written.
- A terminator line **% SZS end BatchProblems**

3.3 Resource Limits

3.3.1 Non-Batch divisions

CPU and wall clock time limits are imposed. The minimal CPU time limit per problem is 240s. The maximal CPU time limit per problem is determined using the relationship used for determining the number of problems, with the minimal number of problems as the *NumberOfProblems*. The CPU time limit is chosen as a reasonable value within the range allowed, and is announced at the competition. The wall clock time limit is imposed in addition to the CPU time limit, to limit very high memory usage that causes swapping. The wall clock time limit per problem is double the CPU time limit. An additional memory limit of 6GB was imposed. The time limits are imposed individually on each solution attempt.

In the demonstration division, each entrant can choose to use either a CPU or a wall clock time limit, whose value is the CPU time limit of the competition divisions.

3.3.2 MZR@Turing division

For each batch there is a there is an overall wall clock time limit, which is available as a command line parameter. The overall limit is the at least 30s multiplied by the number of

problems in the division. There are no CPU time limits.

3.3.3 LTB division

For each batch there is a wall clock time limit per problem, which is provided in the configuration section at the start of each batch. The minimal wall clock time limit per problem is 30s. For each problem category there is an overall wall clock time limit, which is available as a command line parameter. The overall limit is the sum over the batches of the batch's per-problem limit multiplied by the number of problems in the batch. Time spent before starting the first problem of a batch (e.g., preloading and analysing the batch axioms), and times spent between ending a problem and starting the next (e.g., learning from a proof just found), are not part of the times taken on the individual problems, but are part of the overall time taken. There are no CPU time limits.

4 System Evaluation

For each ATP system, for each problem, four items of data are recorded: whether or not the problem was solved, the CPU time taken (not in the MZR@Turing division), the wall clock time taken (not in the MZR@Turing division), and whether or not a solution (proof or model) was output. In the LTB division, time spent before starting the first problem, and times spent between ending a problem and starting the next, are not part of the time taken on problems.

The systems are ranked in the competitions division, from the performance data. The THF, TFA, CNF, EPR, and LTB divisions have an *assurance* ranking class, ranked according to the number of problems solved, but not necessarily accompanied by a proof or model (thus giving only an assurance of the existence of a proof/model). The CASC@Turing, FOF, and FNT divisions have a *proof/model* ranking class, ranked according to the number of problems solved with an acceptable proof/model output. Ties are broken according to the average time over problems solved (CPU time for the non-batch divisions, wall clock time for the batch divisions). In the competition divisions, class winners were announced and prizes are awarded.

- Google provided prize money for CASC@Turing. Prizes were awarded for the proof and model ranking classes. In each case the winner receives £1500, the second place £1000, and the third place £500.
- The Isabelle group at the Technische Universität München provided a travel prize for the ISA problem category of the LTB division. The prize was awarded according to the number of problems solved with a list of axioms sufficient for a proof output. The winner was invited to visit the group at the university for up to one week. The travel and hotel expenses were covered.
- Rearden Commerce provided \$3000 of prize money for the SMO category of the LTB division. Prizes were awarded for the assurance ranking class, and also according to the question answering measure described below. In each case the winner received \$750, the second place \$500, and the third place \$250.

The competition panel decides whether or not the systems' proofs and models are acceptable for the proof/model ranking classes. The criteria include:

- Derivations must be complete, starting at formulae from the problem, and ending at the conjecture (for axiomatic proofs) or a *false* formula (for proofs by contradiction, including CNF refutations).

- For proofs of FOF problems by CNF refutation, the conversion from FOF to CNF must be adequately documented.
- Derivations must show only relevant inference steps.
- Inference steps must document the parent formulae, the inference rule used, and the inferred formula.
- Inference steps must be reasonably fine-grained.
- An unsatisfiable set of ground instances of clauses is acceptable for establishing the unsatisfiability of a set of clauses.
- Models must be complete, documenting the domain, function maps, and predicate maps. The domain, function maps, and predicate maps may be specified by explicit ground lists (of mappings), or by any clear, terminating algorithm.

In the assurance ranking classes the ATP systems are not required to output solutions (proofs or models). However, systems that do output solutions are highlighted in the presentation of results.

In addition to the ranking criteria, other measures are made and presented in the results:

- The *state-of-the-art contribution* (SOTAC) quantifies the unique abilities of each system. For each problem solved by a system, its SOTAC for the problem is the inverse of the number of systems that solved the problem. A system's overall SOTAC is its average SOTAC over the problems it solves.
- The *efficiency measure* balances the number of problems solved with the CPU time taken. It is the average of the inverses of the times for problems solved (CPU times for the non-batch divisions, wall clock times for the LTB division, with times less than the timing granularity rounded up to the granularity, to avoid skewing caused by very low times) multiplied by the fraction of problems solved. This can be interpreted intuitively as the average of the solution rates for problems solved, multiplied by the fraction of problems solved.
- In the LTB division, which uses a wall clock time limit, the *core usage* is the average of the ratios of CPU time to wall clock time used, over the problems solved. This measures the extent to which the systems take advantage the multiple cores.
- In the ISA problem category of the LTB division, the *number of axiom lists* (output of a list of axioms sufficient for a proof) is counted. A well-formed proof can be used to provide such a list.
- In the SMO problem category of the LTB division, the number of *questions answered* (output of the bindings for the outermost existentially quantified variables) is counted (see Section 6.1.3).

At some time after the competition, all high ranking systems in each division are tested over the entire TPTP. This provides a final check for soundness (see Section 6.1 regarding soundness checking before the competition). If a system is found to be unsound during or after the competition, but before the competition report is published, and it cannot be shown that the unsoundness did not manifest itself in the competition, then the system is retrospectively disqualified. At some time after the competition, the proofs and models from the winners of the proof/model ranking classes are checked by the panel. If any of the proofs or models are unacceptable, i.e., they are significantly worse than the samples provided, then that system is retrospectively disqualified. All disqualifications are explained in the competition report.

5 System Entry

To be entered into CASC, systems must be registered using the CASC system registration form. No registrations are accepted after the registration deadline. For each system entered, an entrant has to be nominated to handle all issues (including execution difficulties) arising before and during the competition. The nominated entrant must formally register for CASC. It is not necessary for entrants to physically attend the competition.

Systems can be entered at only the division level, and can be entered into more than one division (a system that is not entered into a competition division is assumed to perform worse than the entered systems, for that type of problem - wimping out is not an option). Entering many similar versions of the same system is deprecated, and entrants may be required to limit the number of system versions that they enter. Systems that rely essentially on running other ATP systems without adding value are deprecated; the competition panel may disallow or move such systems to the demonstration division. The division winners from the previous CASC are automatically entered into their divisions, to provide benchmarks against which progress can be judged.

5.1 System Description

A system description has to be provided for each ATP system entered, using the HTML schema supplied on the CASC web site. (See Section 7 for these descriptions.) The schema has the following sections:

- **Architecture.** This section introduces the ATP system, and describes the calculus and inference rules used.
- **Strategies.** This section describes the search strategies used, why they are effective, and how they are selected for given problems. Any strategy tuning that is based on specific problems' characteristics must be clearly described (and justified in light of the tuning restrictions described in Section 6.1).
- **Implementation.** This section describes the implementation of the ATP system, including the programming language used, important internal data structures, and any special code libraries used. The availability of system is described here.
- **Expected competition performance.** This section makes some predictions about the performance of the ATP system in each of the divisions and categories in which it is competing.
- **References.**

The system description has to be emailed to the competition organizers by the system description deadline. The system descriptions, along with information regarding the competition design and procedures, form the proceedings for the competition.

5.2 Sample Solutions

For systems in the proof/model classes, representative sample solutions must be emailed to the competition organizers by the sample solutions deadline. Use of the TPTP format for proofs and finite interpretations is encouraged. Proof samples for the FOF proof class must include a proof for SEU140+2. Model samples for the FNT model class must include models for NLP042+1 and SWV017+1. The sample solutions must illustrate the use of all inference rules. An explanation must be provided for any non-obvious features.

For systems competing for the ISA problem category prize in the LTB division, representative sample proofs or lists of axioms must be emailed to the competition organizers by the sample solutions deadline. Use of the SZS standards is required. Samples must include a proof or list for SEU140+2. For systems competing for the SMO problem category prize in the LTB division, representative sample answers must be emailed to the competition organizers by the sample solutions deadline. Samples must include an answer for CSR082+1.

6 System Requirements

6.1 System Properties

Entrants must ensure that their systems execute in a competition-like environment, and have the following properties. Entrants are advised to check these properties well in advance of the system delivery deadline. This gives the competition organizers time to help resolve any difficulties encountered. Entrants do not have access to the competition computers.

6.1.1 Soundness and Completeness

- Systems must be sound. At some time before the competition all the systems in the competition divisions are tested for soundness. Non-theorems are submitted to the systems in the FOF@Turing, MZR@Turing, THF, TFA, FOF, EPR, and LTB divisions, and theorems are submitted to the systems in the FNT@Turing, FNT and EPR divisions. Finding a proof of a non-theorem or a disproof of a theorem indicates unsoundness. If a system fails the soundness testing it must be repaired by the unsoundness repair deadline or be withdrawn. The soundness testing eliminates the possibility of a system simply delaying for some amount of time and then claiming to have found a solution. For systems running on entrant supplied computers in the demonstration division, the entrant must perform the soundness testing and report the results to the competition organizers.
- Systems do not have to be complete in any sense, including calculus, search control, implementation, or resource requirements.
- All techniques used must be general purpose, and expected to extend usefully to new unseen problems. The precomputation and storage of information about individual TPTP problems and axiom sets is not allowed. Strategies and strategy selection based on individual TPTP problems is not allowed. If machine learning procedures are used, the learning must ensure that sufficient generalization is obtained so that no there is no specialization to individual problems or their solutions.
 - For the MZR@Turing division, the problems were taken from the publicly available MPTP2078 problem set. As such, precomputation and storage of information about problems in that set, or their solutions, is not directly allowed. However ...
 - Training problems and solutions that were taken from that set could be used for tuning and training systems for the MZR@Turing division. The training problems did not appear in the MZR@Turing division. There were at least twice as many training problems as there were MZR@Turing division problems.
 - The training problems and solutions could be used for producing generally useful strategies that extend to “unseen” problems in the MPTP2078 set. Such strategies can rely on the consistent naming of symbols and formulas in the MPTP2078 set, and could use techniques for memorization and generalization of problems and solutions in the training set. The system description had to fully explain any such tuning or training that was done.

- The competition panel could disqualify any system whose tuning or training was deemed to be problem specific rather than general purpose.
- The system’s performance must be reproducible by running the system again.

6.1.2 Execution

- Systems must run on a single locally provided standard UNIX computer (the *competition computers* - see Section 3.1). ATP systems that cannot run on the competition computers can be entered into the demonstration division.
- Systems must be executable by a single command line, using an absolute path name for the executable, which might not be in the current directory. In the non-batch divisions the command line arguments are the absolute path name of a symbolic link as the problem file name, the time limit (if required by the entrant), and entrant specified system switches. In the batch divisions the command line arguments are the absolute path name of the batch specification file, the overall category time limit (if required by the entrant), and entrant specified system switches. No shell features, such as input or output redirection, may be used in the command line. No assumptions may be made about the format of file names.
- Systems must be fully automatic, i.e., all command line switches have to be the same for all problems in each division.

6.1.3 Output

- In the non-batch divisions all solution output must be to `stdout`. In the batch divisions all solution output must be to the named output file for each problem.
- In the LTB division the systems must print SZS notification lines to `stdout` when starting and ending work on a problem (including any cleanup work, such as deleting temporary files). It is recommended that the result for the problem be output as the last thing before the ending notification line (note, the result must also be output to the solution file anyway). For example

```
% SZS status Started for /home/graph/tptp/TPTP/Problems/CSR/CSR075+2.p
... (system churns away, result and solution output to file)
% SZS status Theorem for /home/graph/tptp/TPTP/Problems/CSR/CSR075+2.p
% SZS status Ended for /home/graph/tptp/TPTP/Problems/CSR/CSR075+2.p
```

- For each problem, the systems must output a distinguished string (specified by the entrant), indicating what solution has been found or that no conclusion has been reached. The distinguished strings should use the SZS ontology and standards [102]. For example

```
SZS status Theorem for SYN075+1
```

or

```
SZS status GaveUp for SYN075+1
```

Regardless of whether the SZS status values are used, the distinguished strings must be different for:

- Proved theorems of FOF problems (SZS status `Theorem`)
- Disproved conjectures of FNT problems (SZS status `CounterSatisfiable`)
- Unsatisfiable sets of formulae (FOF problems without conjectures) and unsatisfiable set of clauses (CNF problems) (SZS status `Unsatisfiable`)

- Satisfiable sets of formulae (FNT problems without conjectures) (SZS status **Satisfiable**)

The first distinguished string output is accepted as the system's result.

- When outputting proofs/models, the start and end of the proof/model must be delimited by distinguished strings (specified by the entrant). The distinguished strings should use the SZS ontology and standards. For example

```
SZS output start CNFRefutation for SYN075-1
...
SZS output end CNFRefutation for SYN075-1
```

Regardless of whether the SZS output forms are used, the distinguished strings must be different for:

- Proofs (SZS output forms **Proof**, **Refutation**, **CNFRefutation**)
- Models (SZS output forms **Model**, **FiniteModel**, **InfiniteModel**, **Saturation**)

The string specifying the problem status must be output before the start of a proof/model. Use of the TPTP format for proofs and finite interpretations is encouraged [112].

- When outputting a list of axioms sufficient for a proof in the ISA problem category of the LTB division, the start and end of the list must be delimited by distinguished strings. The distinguished strings should use the SZS ontology and standards. For example

```
% SZS output start ListOfFOF for SEU104+2
...
% SZS output end ListOfFOF for SEU140+2
```

- When outputting answers for the SMO problem category of the LTB division, the answers must be output using the Tuple or Instantiated answer form of the proposed TPTP standard for answer reporting.

6.1.4 Resource Usage

- The systems that run on the competition computers must be interruptible by a **SIGXCPU** signal, so that the CPU time limit can be imposed, and interruptible by a **SIGALRM** signal, so that the wall clock time limit can be imposed. For systems that create multiple processes, the signal is sent first to the process at the top of the hierarchy, then one second later to all processes in the hierarchy. The default action on receiving these signals is to exit (thus complying with the time limit, as required), but systems may catch the signals and exit of their own accord. If a system runs past a time limit this is noticed in the timing data, and the system is considered to have not solved that problem.
- If an ATP system terminates of its own accord, it may not leave any temporary or intermediate output files. If an ATP system is terminated by a **SIGXCPU** or **SIGALRM**, it may not leave any temporary or intermediate files anywhere other than in **/tmp**. Multiple copies of the ATP systems must be executable concurrently, in the same (NFS cross mounted) directory. It is therefore necessary that temporary files have unique names.
- For practical reasons excessive output from an ATP system is not allowed. A limit, dependent on the disk space available, is imposed on the amount of output that can be produced. The limit is at least 10MB per system.

6.2 System Delivery

For systems running on the competition computers, entrants must email an installation package to the competition organizers by the system delivery deadline. The installation package must be a `.tgz` file containing the system source code, any other files required for installation, and a `ReadMe` file. The `ReadMe` file must contain:

- Instructions for installation
- Instructions for executing the system, using `%s` and `%d` to indicate where the problem file name and time limit must appear in the command line.
- The distinguished strings indicating what solution has been found, and delimiting proofs/models.

The installation procedure may require changing path variables, invoking `make` or something similar, etc, but nothing unreasonably complicated. All system binaries must be created in the installation process; they cannot be delivered as part of the installation package. If the ATP system requires any special software, libraries, etc, which is not part of a standard installation, the competition organizers must be told in the system registration.

The system is installed onto the competition computers by the competition organizers, following the instructions in the `ReadMe` file. Installation failures before the system delivery deadline are passed back to the entrant. (i.e., delivery of the installation package before the system delivery deadline provides an opportunity to fix things if the installation fails!). After the system delivery deadline no further changes or late systems are accepted.

For systems running on entrant supplied computers in the demonstration division, entrants must deliver a source code package to the competition organizers by the start of the competition. The source code package must be a `.tgz` file containing the system source code.

After the competition all competition division systems' source code is made publically available on the CASC web site. In the demonstration division, the entrant specifies whether or not the source code is placed on the CASC web site. An open source license is encouraged.

6.3 System Execution

Execution of the ATP systems on the competition computers is controlled by a `perl` script, provided by the competition organizers. The jobs are queued onto the computers so that each computer is running one job at a time. In the non-LTB divisions, all attempts at the Nth problems in all the divisions and categories are started before any attempts at the (N+1)th problems. In the LTB division all attempts in each category in the division are started before any attempts at the next category.

During the competition a `perl` script parses the systems' outputs. If any of an ATP system's distinguished strings are found then the time used to that point is noted. A system has solved a problem iff it outputs its termination string within the time limit, and a system has produced a proof/model iff it outputs its end-of-proof/model string within the time limit. The result and timing data is used to generate an HTML file, and a web browser is used to display the results.

The execution of the demonstration division systems is supervised by their entrants.

7 The ATP Systems

These system descriptions were written by the entrants.

7.1 CVC4 0.0

Andrew Reynolds¹, Cesare Tinelli¹, Clark Barrett²

¹University of Iowa, USA, ²New York University, USA

Architecture

CVC4 is a Satisfiability Modulo Theories (SMT) solver with support for many theories including linear arithmetic, arrays, bit vectors, and datatypes. Like most SMT solvers, it is based on the DPLL(T) framework [60]. The latest development version of CVC4 supports reasoning for first-order formulas, including a finite model finding feature which can answer satisfiable in the presence of quantified (universal) formulas. Using a theory solver for EUF with cardinality constraints, the system finds a minimal candidate model with respect to uninterpreted sorts, that is, one with a minimal number of ground equivalence classes. When such a candidate model is found, the system will build interpretations for all uninterpreted functions and predicates, using heuristics for defining default values. Consistency is then checked by exhaustively instantiating all quantified formulas with ground instances that may falsify the candidate interpretation. The candidate interpretation is then refined according to these instantiations. This process is repeated until a fixed point is reached, at which point the candidate interpretation specifies a model.

Strategies

The EUF with cardinality constraints theory solver searches for candidate models for a sort S incrementally, i.e. it first postulates that the cardinality of $S = 1$, and if it fails will postulate the cardinality of $S = 2$, etc. Given cardinality constraint k for S , it maintains a disequality graph induced by the ground constraints between terms of sort S . The solver attempts to merge equivalence classes of type S using splitting on demand [9] until no further equivalence classes can be merged. The solver will guide the search by adding conflict lemmas of the form $(C =_{\neq} \text{card}(S)_{\neq} k)$, where C is the explanation of a clique of size $k+1$ in the disequality graph of S . The strategies for model completion and quantifier instantiation are very preliminary at this point.

Implementation

CVC4 is implemented in C++. Much of the code uses context dependent data structures, whose values automatically update when the DPLL(T) search backtracks. The finite model finding module maintains discrimination-tree-like data structures for storing candidate interpretations for functions and predicates, as well as the set of instantiations produced. Since the CVC4 parser does not support the TPTP syntax, the system relies on a translation from TPTP to SMT2.

Expected Competition Performance

The finite model finding feature in CVC4 is still in development. Much can be done to improve the strategy for producing instantiations, including recognizing when certain axioms do not apply. Currently, the bottleneck for CVC4+finite model finding is in processing the large number of instantiations produced. This problem is made worse by the fact that CVC4 currently has no way of removing lemmas once they are added to the system. Preliminary results on TPTP problems are favorable compared to other SMT solvers, which however are not generally equipped to find models for universal formulas.

7.2 E-Darwin 1.5

Björn Pelzer
University Koblenz-Landau, Germany

Architecture

E-Darwin 1.5 [10, 13] is an automated theorem prover for first order clausal logic with equality. It is a modified version of the Darwin prover [10], intended as a testbed for variants of the Model Evolution calculus [14]. Among other things it implements several different approaches [15, 13] to incorporating equality reasoning into Model Evolution. Three principal data structures are used: the context (a set of rewrite literals), the set of constrained clauses, and the set of derived candidates. The prover always selects one candidate, which may be a new clause or a new context literal, and exhaustively computes inferences with this candidate and the context and clause set, moving the results to the candidate set. Afterwards the candidate is inserted into one of the context or the clause set, respectively, and the next candidate is selected. The inferences are superposition-based. Demodulation and various means of redundancy detection are used as well.

Strategies

The uniform search strategy is identical to the one employed in the original Darwin, slightly adapted to account for derived clauses.

Implementation

E-Darwin is implemented in the functional/imperative language OCaml. Darwin's method of storing partial unifiers has been adapted to equations and subterm positions for the superposition inferences in E-Darwin. A combination of perfect and non-perfect discrimination tree indexes is used to store the context and the clauses. The system has been tested on Unix and is available under the GNU Public License from

<http://userpages.uni-koblenz.de/~bpelzer/edarwin>

Expected Competition Performance

After some bugfixes the performance should be slightly better than last year. While the original Darwin performs strongly in EPR, E-Darwin is more of a generalist, less effective in EPR, yet stronger in the other divisions.

7.3 E-KRHyper 1.3

Björn Pelzer
University Koblenz-Landau, Germany

Architecture

E-KRHyper [76] is a theorem proving and model generation system for first-order logic with equality. It is an implementation of the E-hyper tableau calculus [12], which integrates a superposition-based handling of equality [6] into the hyper tableau calculus [11]. The system is an extension of the KRHyper theorem prover [127], which implements the original hyper tableau calculus.^{i/pj}

An E-hyper tableau is a tree whose nodes are labeled with clauses and which is built up by the application of the inference rules of the E-hyper tableau calculus. The calculus rules

are designed such that most of the reasoning is performed using positive unit clauses. Splitting is done without rigid variables. Instead, variables which would be shared between branches are prevented by ground substitutions, which are guessed from the Herbrand universe and constrained by rewrite rules. Redundancy rules allow the detection and removal of clauses that are redundant with respect to a branch. The hyper extension inference from the original hyper tableau calculus is equivalent to a series of E-hyper tableau calculus inference applications. Therefore the implementation of the hyper extension in KRHyper by a variant of semi-naive evaluation [124] is retained in E-KRHyper, where it serves as a shortcut inference for the resolution of non-equational literals.

Strategies

E-KRHyper uses a uniform search strategy for all problems. The E-hyper tableau is generated depth-first, with E-KRHyper always working on a single branch. Refutational completeness and a fair search control are ensured by an iterative deepening strategy with a limit on the maximum term weight of generated clauses. In the LTB division E-KRHyper sequentially tries three axiom selection strategies: an implementation of Krystof Hoder’s SInE algorithm, another incomplete selection based on the CNF representations of the axioms, and finally the complete axiom set.

Implementation

E-KRHyper is implemented in the functional/imperative language OCaml. The system accepts input in the TPTP-format and in the TPTP-supported Protein-format. The calculus implemented by E-KRHyper works on clauses, so first order formula input is converted into CNF by an algorithm similar to the one used by Otter [58], with some additional connector literals to prevent explosive clause growth when dealing with DNF-like structures. E-KRHyper operates on an E-hyper tableau which is represented by linked node records. Several layered discrimination-tree based indexes (both perfect and non-perfect) provide access to the clauses in the tableau and support backtracking. The system runs on Unix and MS-Windows platforms and is available under the GNU Public License from

<http://www.uni-koblenz.de/~bpelzer/ekrhyper>

Expected Competition Performance

The redundancy handling has been improved, resulting in a slight improvement over last year. Overall E-KRHyper will remain in the middle ground.

7.4 E-MaLeS 1.1

Daniel Kuehlwein¹, Josef Urban¹, Stephan Schulz²

¹Radboud University Nijmegen, The Netherlands, ²Technische Universität München, Germany

Architecture

E-MaLeS (Machine Learning of Strategies) uses kernel-based machine learning to improve the strategy selection of E prover.

From the known performance of different E strategies over the TPTP library, the system learns which strategy is most likely to solve a problem. Given a new problem, E-MaLeS extracts the problems features (number of terms, number of literals, etc) and applies the learned function to the features. The result is a (hopefully optimal) strategy. As a new addition in 1.1, E-MaLeS

learns not only the optimal strategy, but also for how long it should run each strategy. This allows better strategy scheduling.

Strategies

Each E strategy is run for the predicted time, starting with the strategy with the lowest predicted time.

Implementation

E-MaLeS is implemented in python, using the numpy and scipy libraries. It is available from

<http://www.cs.ru.nl/~kuehlwein/>

Expected Competition Performance

E-MaLeS should perform slightly better than E.

7.5 EP 1.4pre

Stephan Schulz

Technische Universität München, Germany

Architecture

E 1.4pre [82, 83] is described in this section. E is a purely equational theorem prover for full first-order logic with equality. It consists of an (optional) clausifier for pre-processing full first-order formulae into clausal form, and a saturation algorithm implementing an instance of the superposition calculus with negative literal selection and a number of redundancy elimination techniques. E is based on the DISCOUNT-loop variant of the *given-clause* algorithm, i.e., a strict separation of active and passive facts. No special rules for non-equational literals have been implemented. Resolution is effectively simulated by paramodulation and equality resolution.

EP 1.4pre is just a combination of E 1.4pre in verbose mode and a proof analysis tool extracting the used inference steps. For the LTB division, a control program uses a SInE-like analysis to extract reduced axiomatizations that are handed to several instances of E.

Strategies

Proof search in E is primarily controlled by a literal selection strategy, a clause evaluation heuristic, and a simplification ordering. The prover supports a large number of pre-programmed literal selection strategies. Clause evaluation heuristics can be constructed on the fly by combining various parametrized primitive evaluation functions, or can be selected from a set of predefined heuristics. Clause evaluation heuristics are based on symbol-counting, but also take other clause properties into account. In particular, the search can prefer clauses from the set of support, or containing many symbols also present in the goal. Supported term orderings are several parametrized instances of Knuth-Bendix-Ordering (KBO) and Lexicographic Path Ordering (LPO).

The automatic mode is based on a static partition of the set of all clausal problems based on a number of simple features (number of clauses, maximal symbol arity, presence of equality, presence of non-unit and non-Horn clauses,...). Each class of clauses is automatically assigned a heuristic that performs well on problems from this class in test runs. About 100 different strategies have been evaluated on all untyped first-order problems from TPTP 4.1.0.

Implementation

E is build around perfectly shared terms, i.e. each distinct term is only represented once in a term bank. The whole set of terms thus consists of a number of interconnected directed acyclic graphs. Term memory is managed by a simple mark-and-sweep garbage collector. Unconditional (forward) rewriting using unit clauses is implemented using perfect discrimination trees with size and age constraints. Whenever a possible simplification is detected, it is added as a rewrite link in the term bank. As a result, not only terms, but also rewrite steps are shared. Subsumption and contextual literal cutting (also known as subsumption resolution) is supported using feature vector indexing [83]. Superposition and backward rewriting use fingerprint indexing, a new technique combining ideas from feature vector indexing and path indexing. Finally, LPO and KBO are implemented using the elegant and efficient algorithms developed by Bernd Löchner in [55, 55]. The prover and additional information are available at

<http://www.eprover.org>

Expected Competition Performance

EP 1.4pre is the CASC-23 CNF division winner.

7.6 EP 1.6pre

Stephan Schulz
Technische Universität München, Germany

Architecture

E 1.6pre [82] is a purely equational theorem prover for full first-order logic with equality. It consists of an (optional) clausifier for pre-processing full first-order formulae into clausal form, and a saturation algorithm implementing an instance of the superposition calculus with negative literal selection and a number of redundancy elimination techniques. E is based on the DISCOUNT-loop variant of the *given-clause* algorithm, i.e., a strict separation of active and passive facts. No special rules for non-equational literals have been implemented. Resolution is effectively simulated by paramodulation and equality resolution.

EP 1.6pre is just a combination of E 1.6pre in verbose mode and a proof analysis tool extracting the used inference steps. For the LTB and MZR divisions, a control program uses a SInE-like analysis to extract reduced axiomatizations that are handed to several instances of E.

Strategies

Proof search in E is primarily controlled by a literal selection strategy, a clause evaluation heuristic, and a simplification ordering. The prover supports a large number of pre-programmed literal selection strategies. Clause evaluation heuristics can be constructed on the fly by combining various parametrized primitive evaluation functions, or can be selected from a set of predefined heuristics. Clause evaluation heuristics are based on symbol-counting, but also take other clause properties into account. In particular, the search can prefer clauses from the set of support, or containing many symbols also present in the goal. Supported term orderings are several parametrized instances of Knuth-Bendix-Ordering (KBO) and Lexicographic Path Ordering (LPO).

The automatic mode is based on a static partition of the set of all clausal problems based on a number of simple features (number of clauses, maximal symbol arity, presence of equality,

presence of non-unit and non-Horn clauses,...). Each class of clauses is automatically assigned a heuristic that performs well on problems from this class in test runs. About 60 different strategies have been evaluated on all bug-free untyped first-order problems from TPTP 5.3.0.

Implementation

E is build around perfectly shared terms, i.e. each distinct term is only represented once in a term bank. The whole set of terms thus consists of a number of interconnected directed acyclic graphs. Term memory is managed by a simple mark-and-sweep garbage collector. Unconditional (forward) rewriting using unit clauses is implemented using perfect discrimination trees with size and age constraints. Whenever a possible simplification is detected, it is added as a rewrite link in the term bank. As a result, not only terms, but also rewrite steps are shared. Subsumption and contextual literal cutting (also known as subsumption resolution) is supported using feature vector indexing [83]. Superposition and backward rewriting use fingerprint indexing [84], a new technique combining ideas from feature vector indexing and path indexing. Finally, LPO and KBO are implemented using the elegant and efficient algorithms developed by Bernd Löchner in [55, 55]. The prover and additional information are available at

<http://www.eprover.org>

Expected Competition Performance

E 1.6pre has improved indexing and now integrates optional SinE-like preprocessing for large FOF problems. The system is expected to perform well in most proof classes, but will at best complement top systems in the disproof classes.

7.7 FIMO 0.3

Orkunt Sabuncu

University of Potsdam, Germany

Architecture

Fimo is a system for computing finite models of first-order formulas by incremental Answer Set Programming (iASP). The input theory is transformed to an incremental logic program. If any, answer sets of this program represent finite models of the input theory. iClingo is used for computing answer sets of iASP programs.

Strategies

Fimo is the successor of the system `fmc2iasp` [43]. Unlike `fmc2iasp` Fimo does not rely on flattening for translating the input theory to a satisfiability problem.

Fimo features symmetry breaking and incremental answer set solving provided by the underlying iASP system iClingo. Additionally, Fimo deskolemizes some of the Skolem functions by transforming them to aggregates in the resulting logic program.

Implementation

Fimo is developed in Python. It is available from

<http://potassco.sourceforge.net>

Expected Competition Performance

We are expecting a similar performance to the version that competed in CASC 2011.

7.8 iProver 0.9

Konstantin Korovin
University of Manchester, United Kingdom

Architecture

iProver is an automated theorem prover based on an instantiation calculus Inst-Gen [41, 51] which is complete for first-order logic. One of the distinctive features of iProver is a modular combination of first-order reasoning with ground reasoning. In particular, iProver currently integrates MiniSat [40] for reasoning with ground abstractions of first-order clauses. In addition to instantiation, iProver implements ordered resolution calculus and a combination of instantiation and ordered resolution; see [50] for the implementation details. The saturation process is implemented as a modification of a given clause algorithm. iProver uses non-perfect discrimination trees for the unification indexes, priority queues for passive clauses, and a compressed vector index for subsumption and subsumption resolution (both forward and backward). The following redundancy eliminations are implemented: blocking non-proper instantiations; mismatching constraints [42, 50]; global subsumption [50]; resolution-based simplifications and propositional-based simplifications. A compressed feature vector index is used for efficient forward/backward subsumption and subsumption resolution. Equality is dealt with (internally) by adding the necessary axioms of equality with an option of using Brand's transformation. In the LTB division, iProver-SInE uses axiom selection based on the SInE algorithm [47] as implemented in Vampire [45], i.e., axiom selection is done by Vampire and proof attempts are done by iProver.

Major additions in the current version are:

- answer computation,
- several modes for model output using first-order definitions in term algebra,
- Brand's transformation.

Strategies

iProver has around 40 options to control the proof search including options for literal selection, passive clause selection, frequency of calling the SAT solver, simplifications and options for combination of instantiation with resolution. At CASC iProver will execute a small number of fixed schedules of selected options depending on general syntactic properties such as Horn/non-Horn, equational/non-equational, and maximal term depth.

Implementation

iProver is implemented in OCaml and for the ground reasoning uses MiniSat. iProver accepts FOF and CNF formats, where Vampire [45] is used for clausification of FOF problems.

iProver is available from

<http://www.cs.man.ac.uk/~korovink/iprover/>

Expected Competition Performance

iProver 0.9 is the CASC-23 EPR division winner.

7.9 iProver 0.99

Konstantin Korovin, Christoph Stickel
University of Manchester, United Kingdom

Architecture

iProver is an automated theorem prover based on an instantiation calculus Inst-Gen [41, 51] which is complete for first-order logic. iProver combines first-order reasoning with ground reasoning for which it uses MiniSat [40]. iProver also combines instantiation with ordered resolution; see [50] for the implementation details. The proof search is implemented using a saturation process based on the given clause algorithm. iProver uses non-perfect discrimination trees for the unification indexes, priority queues for passive clauses, and a compressed vector index for subsumption and subsumption resolution (both forward and backward). The following redundancy eliminations are implemented: blocking non-proper instantiations; dismatching constraints [GK04,Kor08]; global subsumption [50]; resolution-based simplifications and propositional-based simplifications. A compressed feature vector index is used for efficient forward/backward subsumption and subsumption resolution. Equality is dealt with (internally) by adding the necessary axioms of equality with an option of using Brand's transformation.

In the LTB division, iProver uses axiom selection based on the SInE algorithm [47] as implemented in Vampire [46], i.e., axiom selection is done by Vampire and proof attempts are done by iProver.

iProver features are summaries below with recent additions marked with (*).

- proof extraction for both instantiation and resolution (*),
- model representation, using first-order definitions in term algebra,
- answer substitutions,
- semantic filtering (*),
- type inference [36] (*),
- Brand's transformation.

Type inference is targeted at improving finite model finding and symmetry breaking. Semantic filtering is used in preprocessing to eliminate irrelevant clauses. Proof extraction is challenging due to simplifications such as global subsumption which involve global reasoning with the whole clause set and can be computationally expensive.

Strategies

iProver has around 60 options to control the proof search including options for literal selection, passive clause selection, frequency of calling the SAT solver, simplifications and options for combination of instantiation with resolution. At CASC iProver will execute a small number of fixed schedules of selected options depending on general syntactic properties such as Horn/non-Horn, equational/non-equational, and maximal term depth. The strategy for satisfiable problems (FNT division) includes finite model finding.

Implementation

iProver is implemented in OCaml and for the ground reasoning uses MiniSat [40]. iProver accepts FOF and CNF formats. Vampire [46] is used for proof-producing clausification of FOF problems as well as for axiom selection [47] in the LTB division. iProver is available from

<http://www.cs.man.ac.uk/~korovink/iprover/>

Expected Competition Performance

We expect very good performance on satisfiable problems (FNT division) due to recent additions. In the EPR division we expect similar performance to the previous version but additional computations needed for the proof extraction can degrade the performance.

7.10 iProver-Eq 0.8

Christoph Stickse, Konstantin Korovin
University of Manchester, United Kingdom

Architecture

iProver-Eq [52] extends the iProver system [50] with built-in equational reasoning, along the lines of [42]. As in the iProver system, first-order reasoning is combined with ground satisfiability checking where the latter is delegated to an off-the-shelf ground solver.

iProver-Eq consists of three core components: i) ground reasoning by an SMT solver, ii) first-order equational reasoning on literals in a candidate model by a labelled unit superposition calculus [52, 52] and iii) instantiation of clauses with substitutions obtained by ii).

Given a set of first-order clauses, iProver-Eq first abstracts it to a set of ground clauses which are then passed to the ground solver. If the ground abstraction is unsatisfiable, then the set of first-order clauses is also unsatisfiable. Otherwise, literals are selected from the first-order clauses based on the model of the ground solver. The labelled unit superposition calculus checks whether selected literals are conflicting. If they are conflicting, then clauses are instantiated such that the ground solver has to refine its model in order to resolve the conflict. Otherwise, satisfiability of the initial first-order clause set is shown.

Clause selection and literal selection in the unit superposition calculus are implemented in separate given clause algorithms. Relevant substitutions are accumulated in labels during unit superposition inferences and then used to instantiate clauses. For redundancy elimination iProver-Eq uses demodulation, dismatching constraints and global subsumption. In order to efficiently propagate redundancy elimination from instantiation into unit superposition, we implemented different representations of labels based on sets, AND/OR-trees and OBDDs. Non-equational resolution and equational superposition inferences provide further simplifications.

For the LTB division, iProver-Eq-SInE runs iProver-Eq as the underlying inference engine of SInE [47], i.e., axiom selection is done by SInE, and proof attempts are done by iProver-Eq.

Strategies

Proof search options in iProver-Eq control clause and literal selection in the respective given clause algorithms. Equally important is the global distribution of time between the inference engines and the ground solver. At CASC, iProver-Eq will execute a fixed schedule of selected options.

If no equational literals occur in the input, iProver-Eq falls back to the inference rules of iProver, otherwise the latter are disabled and only unit superposition is used. If all clauses are unit equations, no instances need to be generated and the calculus is run without the otherwise necessary bookkeeping.

Implementation

iProver-Eq is implemented in OCaml and uses a prototype of the CVC4 SMT solver, which is the successor of CVC3 [8], for the ground reasoning in the equational case and MiniSat [40]

in the non-equational case. iProver-Eq accepts FOF and CNF formats, where Vampire [46] is used for clausification and preprocessing of FOF problems.

iProver-Eq is available from

<http://www.cs.man.ac.uk/~sticksec/iprover-eq>

Expected Competition Performance

iProver-Eq has seen many optimisations from the version in the previous CASCs, in particular regarding labelling and the integration of CVC4. We expect reasonably good performance in all divisions, including the EPR divisions where instantiation-based methods are particularly strong.

7.11 Isabelle 2012

Jasmin C. Blanchette¹, Lawrence C. Paulson², Tobias Nipkow¹, Makarius Wenzel³

¹Technische Universität München, Germany, ²University of Cambridge, United Kingdom,

³Université Paris Sud, France

Architecture

Isabelle/HOL 2012 [62] is the higher-order logic incarnation of the generic proof assistant Isabelle2012. Isabelle/HOL provides several automatic proof tactics, notably an equational reasoner [61], a classical reasoner [75], and a tableau prover [73]. It also integrates external first- and higher-order provers via its subsystem Sledgehammer [74, 24].

Previous versions of Isabelle relied on the TPTP2X tool to translate TPTP files to Isabelle theory files. Starting this year, Isabelle includes a parser for the TPTP syntaxes CNF, FOF, TFF0, and THF0 as well as TPTP versions of its popular tools, invocable on the command line as `isabelle tptp_tool max_secs file.p`. For example:

```
isabelle tptp\_isabelle\_demo 100 SEU/SEU824^3.p
```

Two versions of Isabelle participate this year. The *demo* version includes its competitors LEO-II [18] and Satallax [33] as Sledgehammer backends, whereas the *competition* version leaves them out.

Strategies

The *Isabelle* tactic submitted to the competition simply tries the following tactics sequentially:

- `sledgehammer` – Invokes the following sequence of provers as oracles via Sledgehammer:
 - `satallax` – Satallax 2.4 [33] (*demo only*);
 - `leo2` – LEO-II 1.3.2 [18] (*demo only*);
 - `spass` – SPASS 3.8ds [25];
 - `vampire` – Vampire 1.8 (revision 1435) [79];
 - `e` – E 1.4 [83];
 - `z3_tptp` – Z3 3.2 with TPTP syntax [38].
- `nitpick` – For problems involving only the type `$o` of Booleans, checks whether a finite model exists using Nitpick [27].
- `simp` – Performs equational reasoning using rewrite rules [61].

- **blast** – Searches for a proof using a fast untyped tableau prover and then attempts to reconstruct the proof using Isabelle tactics [73].
- **auto+spass** – Combines simplification and classical reasoning [75] under one roof; then invoke Sledgehammer with SPASS on any subgoals that emerge.
- **z3** – Invokes the SMT solver Z3 3.2 [38].
- **cvc3** – Invokes the SMT solver CVC3 2.2 [8].
- **fast** – Searches for a proof using sequent-style reasoning, performing a depth-first search [75]. Unlike **blast**, it constructs proofs directly in Isabelle. That makes it slower but enables it to work in the presence of the more unusual features of HOL, such as type classes and function unknowns.
- **best** – Similar to **fast**, except that it performs a best-first search.
- **force** – Similar to **auto**, but more exhaustive.
- **meson** – Implements Loveland’s MESON procedure [56]. Constructs proofs directly in Isabelle.
- **fastforce** – Combines **fast** and **force**.

Implementation

Isabelle is a generic theorem prover written in Standard ML. Its meta-logic, Isabelle/Pure, provides an intuitionistic fragment of higher-order logic. The HOL object logic extends pure with a more elaborate version of higher-order logic, complete with the familiar connectives and quantifiers. Other object logics are available, notably FOL (first-order logic) and ZF (Zermelo-Fraenkel set theory).

The implementation of Isabelle relies on a small LCF-style kernel, meaning that inferences are implemented as operations on an abstract **theorem** datatype. Assuming the kernel is correct, all values of type **theorem** are correct by construction.

Most of the code for Isabelle was written by the Isabelle teams at the University of Cambridge and the Technische Universität München. Isabelle/HOL is available for all major platforms under a BSD-style license from

<http://www.cl.cam.ac.uk/research/hvg/Isabelle>

Expected Competition Performance

By integrating LEO-II and Satallax as two of many of its tactics, Isabelle should have all the tools at its disposal to win the competition this year.

7.12 leanCoP 2.2

Jens Otten
University of Potsdam, Germany

Architecture

leanCoP [65, 63] is an automated theorem prover for classical first-order logic with equality. It is a very compact implementation of the connection (tableau) calculus [21, 54].

Strategies

The reduction rule of the connection calculus is applied before the extension rule. Open branches are selected in a depth-first way. Iterative deepening on the proof depth is used to achieve completeness. Additional inference rules and strategies include regularity, lemmata, and restricted

backtracking [64]. leanCoP uses an optimized structure-preserving transformation into clausal form [64] and a fixed strategy scheduling, which is invoked by a shell script.

Implementation

leanCoP is implemented in Prolog. The source code of the core prover is only a few lines long and fits on half a page. Prolog's built-in indexing mechanism is used to quickly find connections.

leanCoP can read formulae in leanCoP syntax as well as in (raw) TPTP first-order syntax. Equality axioms and axioms to support distinct objects are automatically added if required. The leanCoP core prover returns a very compact connection proof, which can be translated into different proof formats, e.g., into a lean (unofficial) TPTP syntax format for representing connection proofs [66] or into a readable text proof.

The leanCoP core prover and all its additional components run on any (standard) Prolog system; ECLiPSe, SICStus and SWI Prolog are currently explicitly supported. The shell script that implements the fixed strategy scheduling runs on Linux/Unix and MacOS platforms as well as on most Windows platforms.

The source code of leanCoP 2.2 is available under the GNU general public license. Together with more information it can be found on the leanCoP website at

<http://www.leancop.de>

The website also contains information about the leanCoP versions ileanCoP and MleanCoP, leading automated theorem provers for first-order intuitionistic logic and several first-order modal logics, respectively.

Expected Competition Performance

As the core prover has not changed, the performance of leanCoP 2.2 is expected to be similar to the performance of leanCoP 2.1.

7.13 leanCoP-ARDE 2.2

Mario Frank, Thomas Raths, Jens Otten
University of Potsdam, Germany

Architecture

ARDE (Axiom Relevance Decision Engine, current version 0.5) is a tool that processes large sets of axioms and categorizes them into classes. Given a conjecture (problem file) it selects those axioms from the axiom set that are likely required for a proof of the conjecture. For the proof process itself it consults leanCoP [65, 63], a compact theorem prover for first-order classical logic.

Strategies

ARDE analyses all axioms and gathers information about included predicate symbols, their arity and additional information like their polarities. Based on this information, it tries to decide which axioms are likely needed for a proof of the conjecture. These axioms together with the conjecture is then given to the leanCoP 2.2 theorem prover.

Implementation

ARDE is implemented in C++ using the C++11 standard. It uses version 1.46 of the Boost library (e.g. Spirit) and GCC 4.6.3. Other used tools include egrep and GNU make.

Expected Competition Performance

On large problems containing many axioms leanCoP-ARDE 2.2 is expected to show a better performance than the leanCoP 2.2 prover on its own.

7.14 LEO-II 1.4

Christoph Benzmüller
Free University Berlin, Germany

Architecture

LEO-II [18], the successor of LEO [17], is a higher-order ATP system based on extensional higher-order resolution. More precisely, LEO-II employs a refinement of extensional higher-order RUE resolution [16]. LEO-II is designed to cooperate with specialist systems for fragments of higher-order logic. By default, LEO-II cooperates with the first-order ATP system E [82]. LEO-II is often too weak to find a refutation amongst the steadily growing set of clauses on its own. However, some of the clauses in LEO-II's search space attain a special status: they are first-order clauses modulo the application of an appropriate transformation function. Therefore, LEO-II launches a cooperating first-order ATP system every n iterations of its (standard) resolution proof search loop (e.g., 10). If the first-order ATP system finds a refutation, it communicates its success to LEO-II in the standard SZS format. Communication between LEO-II and the cooperating first-order ATP system uses the TPTP language and standards.

Strategies

LEO-II employs an adapted "Otter loop". Moreover, LEO-II now also uses some very basic strategy scheduling to try different search strategies or flag settings. These search strategies also include some different relevance filters.

Implementation

LEO-II is implemented in Objective Caml version 3.12, and its problem representation language is the TPTP THF language [19]. In fact, the development of LEO-II has largely paralleled the development of the TPTP THF language and related infrastructure [111].

LEO-II's parser supports the TPTP THF0 language and also the TPTP languages FOF and CNF.

Unfortunately the LEO-II system still uses only a very simple sequential collaboration model with first-order ATPs instead of using the more advanced, concurrent and resource-adaptive OANTS architecture [20] as exploited by its predecessor LEO.

The LEO-II system is distributed under a BSD style license, and it is available from

<http://www.leoprover.org>

Expected Competition Performance

LEO-II has not improved much since 2010. The main modifications concern proof output in order to enable proof reconstruction/verification of LEO-II proofs in other systems [85].

7.15 MaLAREa 0.4

Josef Urban¹, Daniel Kuehlwein¹, Stephan Schulz², Jiri Vyskocil³

¹Radboud University Nijmegen, The Netherlands, ²Technische Universität München, Germany,

³Czech Technical University, Czech Republic

Architecture

MaLAREa 0.4 [125, 126] is a metasystem for ATP in large theories where symbol and formula names are used consistently. It uses several deductive systems (now E, SPASS, Vampire, Paradox, Mace), as well as complementary AI techniques like machine learning (the SNoW system) based on symbol-based similarity, model-based similarity, term-based similarity, and obviously previous successful proofs. The version for CASC 2012 will mainly use E prover, possibly also Prover9, Mace and Paradox. The premise selection methods will likely also use kernel-based learning and E's implementation of SInE.

Strategies

The basic strategy is to run ATPs on problems, then use the machine learner to learn axiom relevance for conjectures from solutions, and use the most relevant axioms for next ATP attempts. This is iterated, using different timelimits and axiom limits. Various features are used for learning, and the learning is complemented by other criteria like model-based reasoning, symbol and term-based similarity, etc.

Implementation

The metasystem is implemented in ca. 2500 lines of Perl. It uses many external programs - the above mentioned ATPs and machine learner, TPTP utilities, LADR utilities for work with models, and some standard Unix tools.

MaLAREa is available from

<https://github.com/JUrban/MPTP2/tree/master/MaLAREa>

The metasystem's Perl code is released under GPL2.

Expected Competition Performance

Thanks to machine learning, MaLAREa is strongest on batches of many related problems with many redundant axioms where some of the problems are easy to solve and can be used for learning the axiom relevance. MaLAREa is not very good when all problems are too difficult (nothing to learn from), or the problems (are few and) have nothing in common. Some of its techniques (selection by symbol and term-based similarity, model-based reasoning) could however make it even there slightly stronger than standard ATPs. MaLAREa has a very good performance on the MPTP Challenge, which is a predecessor of the LTB division, and it is the winner of the 2008 MZR LTB category which had the same rules as the MPTP Challenge. Since then the LTB rules changed and MaLAREa did not compete.

7.16 Muscadet 4.2

Dominique Pastre
University Paris Descartes, France

Architecture

The Muscadet theorem prover is a knowledge-based system. It is based on Natural Deduction, following the terminology of [26] and [67], and uses methods which resembles those used by humans. It is composed of an inference engine, which interprets and executes rules, and of one or several bases of facts, which are the internal representation of "theorems to be proved". Rules are either universal and put into the system, or built by the system itself by metarules from data (definitions and lemmas). Rules may add new hypotheses, modify the conclusion, create objects, split theorems into two or more subtheorems or build new rules which are local for a (sub-)theorem.

Strategies

There are specific strategies for existential, universal, conjunctive or disjunctive hypotheses and conclusions, and equalities. Functional symbols may be used, but an automatic creation of intermediate objects allows deep subformulae to be flattened and treated as if the concepts were defined by predicate symbols. The successive steps of a proof may be forward deduction (deduce new hypotheses from old ones), backward deduction (replace the conclusion by a new one), refutation (only if the conclusion is a negation), search for objects satisfying the conclusion or dynamic building of new rules.

The system is also able to work with second order statements. It may also receive knowledge and know-how for a specific domain from a human user; see [68] and [69]. These two possibilities are not used while working with the TPTP Library.

Implementation

Muscadet [70] is implemented in SWI-Prolog. Rules are written as more or less declarative Prolog clauses. Metarules are written as sets of Prolog clauses. The inference engine includes the Prolog interpreter and some procedural Prolog clauses. A theorem may be split into several subtheorems, structured as a tree with "and" and "or" nodes. All the proof search steps are memorized as facts including all the elements which will be necessary to extract later the useful steps (the name of the executed action or applied rule, the new facts added or rule dynamically built, the antecedents and a brief explanation).

Muscadet is available from

<http://www.math-info.univ-paris5.fr/~pastre/muscadet/muscadet.html>

Expected Competition Performance

The best performances of Muscadet will be for problems manipulating many concepts in which all statements (conjectures, definitions, axioms) are expressed in a manner similar to the practice of humans, especially of mathematicians [71, 72]. It will have poor performances for problems using few concepts but large and deep formulas leading to many splittings. Its best results will be in set theory, especially for functions and relations. Its originality is that proofs are given in natural style. This 4.2 version is a minor revision of 4.1 version. It differs only in that a number of bugs have been fixed.

7.17 Nitrox 2012

Jasmin C. Blanchette¹, Emina Torlak²

¹Technische Universität München, Germany, ²University of California, USA

Architecture

Nitrox is the first-order version of Nitpick [27], an open source counterexample generator for Isabelle/HOL [62]. It builds on Kodkod [123], a highly optimized first-order relational model finder based on SAT. The name Nitrox is a portmanteau of *Nitpick* and *Paradox* (clever, eh?).

Strategies

Nitrox employs Kodkod to find a finite model of the negated conjecture. It performs a few transformations on the input, such as pushing quantifiers inside, but 99solver.

The translation from HOL to Kodkod’s first-order relational logic (FORL) is parameterized by the cardinalities of the atomic types occurring in it. Nitrox enumerates the possible cardinalities for the universe. If a formula has a finite counterexample, the tool eventually finds it, unless it runs out of resources.

Nitpick is optimized to work with higher-order logic (HOL) and its definitional principles (e.g., (co)inductive predicates, (co)inductive datatypes, (co)recursive functions). When invoked on untyped first-order problem, few of its optimizations come into play, and the problem handed to Kodkod is essentially a first-order relational logic (FORL) rendering of the TPTP FOF problem. There are two main exceptions:

- Nested quantifiers are moved as far inside the formula as possible before Kodkod gets a chance to look at them [27].
- Definitions invoked with fixed arguments are specialized.

Implementation

Nitrox, like most of Isabelle/HOL, is written in Standard ML. Unlike Isabelle itself, which adheres to the LCF small-kernel discipline, Nitrox does not certify its results and must be trusted. Kodkod is written in Java. MiniSat 1.14 is used as the SAT solver.

Expected Competition Performance

Since Nitpick was designed for HOL, it doesn’t have any type inference à la Paradox. It also doesn’t use the SAT solver incrementally, which penalizes it a bit (but not as much as the missing type inference). Kodkod itself is known to perform less well on FOF than Paradox, because it is designed and optimized for a somewhat different logic, FORL. On the other hand, Kodkod’s symmetry breaking seems better calibrated than Paradox’s. Hence, we expect Nitrox to end up in second place at best in the TNF category.

7.18 Paradox 3.0

Koen Claessen, Niklas Sörensson
Chalmers University of Technology, Sweden

Architecture

Paradox [37] is a finite-domain model generator. It is based on a MACE-style [58] flattening

and instantiating of the first-order clauses into propositional clauses, and then the use of a SAT solver to solve the resulting problem.

Paradox incorporates the following features: Polynomial-time *clause splitting heuristics*, the use of *incremental SAT*, *static symmetry reduction* techniques, and the use of *sort inference*.

Strategies

There is only one strategy in Paradox:

1. Analyze the problem, finding an upper bound N on the domain size of models, where N is possibly infinite. A finite such upper bound can be found, for example, for EPR problems.
2. Flatten the problem, and split clauses and simplify as much as possible.
3. Instantiate the problem for domain sizes 1 up to N , applying the SAT solver incrementally for each size. Report “SATISFIABLE” when a model is found.
4. When no model of sizes smaller or equal to N is found, report “CONTRADICTION”.

In this way, Paradox can be used both as a model finder and as an EPR solver.

Implementation

The main part of Paradox is implemented in Haskell using the GHC compiler. Paradox also has a built-in incremental SAT solver which is written in C++. The two parts are linked together on the object level using Haskell’s Foreign Function Interface.

Expected Competition Performance

Paradox 3.0 is the CASC-23 FNT division winner.

7.19 Princess 2012-05-28

Philipp Rümmer, Aleksandar Zeljic
Uppsala University, Sweden

Architecture

Princess [80, 81] is a theorem prover for first-order logic modulo linear integer arithmetic. The prover has been under development since 2007, and represents a combination of techniques from the areas of first-order reasoning and SMT solving. The main underlying calculus is a free-variable tableau calculus, which is extended with constraints to enable backtracking-free proof expansion, and positive unit hyper-resolution for lightweight instantiation of quantified formulae. Linear integer arithmetic is handled using a set of built-in proof rules resembling the Omega test, which altogether yields a calculus that is complete for full Presburger arithmetic, for first-order logic, and for a number of further fragments.

The internal calculus of Princess only supports uninterpreted predicates; uninterpreted functions are encoded as predicates, together with the usual axioms. Through appropriate translation of quantified formulae with functions, the e-matching technique common in SMT solvers can be simulated; triggers in quantified formulae are chosen based on heuristics similar to those in the Simplify prover.

Strategies

Princess supports a number of different proof expansion strategies (e.g., depth-first, breadth-first), which are chosen based on syntactic properties of a problem (in particular, the kind

of quantifiers occurring in the problem). Further options exist to control, for instance, the selection of triggers in quantified formulae, clausification, and the handling of functions.

For CASC, Princess will run a schedule with a small number of configurations for each problem (portfolio method). The schedule is determined either statically, or dynamically using syntactic attributes of problems (such as number and kind of quantifiers, etc), based on training using a random sample of problems from the TPTP library.

Implementation

Princess is entirely written in Scala and runs on any recent Java virtual machine; besides the standard Scala and Java libraries, only the Cup parser library is employed. Princess is available from

<http://www.philipp.ruemmer.org/princess.shtml>

Expected Competition Performance

Since Princess enters CASC for the first time, performance predictions are entirely speculative. Princess is mainly designed for the TFI category, and should perform reasonably well here. Support for rationals and reals is preliminary and incomplete, so that mediocre results are expected for TFR. Finally, although Princess is in theory complete for FOL, it is not designed for this logic and enters FOF just for fun.

7.20 Prover9 2009-11A

Bob Veroff on behalf of William McCune
University of New Mexico, USA

Architecture

Prover9, Version 2009-11A, is a resolution/paramodulation prover for first-order logic with equality. Its overall architecture is very similar to that of Otter-3.3 [58]. It uses the “given clause algorithm”, in which not-yet-given clauses are available for rewriting and for other inference operations (sometimes called the “Otter loop”).

Prover9 has available positive ordered (and nonordered) resolution and paramodulation, negative ordered (and nonordered) resolution, factoring, positive and negative hyperresolution, UR-resolution, and demodulation (term rewriting). Terms can be ordered with LPO, RPO, or KBO. Selection of the “given clause” is by an age-weight ratio.

Proofs can be given at two levels of detail: (1) standard, in which each line of the proof is a stored clause with detailed justification, and (2) expanded, with a separate line for each operation. When FOF problems are input, proof of transformation to clauses is not given.

Completeness is not guaranteed, so termination does not indicate satisfiability.

Strategies

Like Otter, Prover9 has available many strategies; the following statements apply to CASC-2012.

Given a problem, Prover9 adjusts its inference rules and strategy according to syntactic properties of the input clauses such as the presence of equality and non-Horn clauses. Prover9 also does some preprocessing, for example, to eliminate predicates.

In previous CASC competitions, Prover9 has used LPO to order terms for demodulation and for the inference rules, with a simple rule for determining symbol precedence. For CASC 2012, we are going to use KBO instead.

For the FOF problems, a preprocessing step attempts to reduce the problem to independent subproblems by a miniscope transformation; if the problem reduction succeeds, each subproblem is classified and given to the ordinary search procedure; if the problem reduction fails, the original problem is classified and given to the search procedure.

Implementation

Prover9 is coded in C, and it uses the LADR libraries. Some of the code descended from EQP [57]. (LADR has some AC functions, but Prover9 does not use them). Term data structures are not shared (as they are in Otter). Term indexing is used extensively, with discrimination tree indexing for finding rewrite rules and subsuming units, FPA/Path indexing for finding subsumed units, rewritable terms, and resolvable literals. Feature vector indexing [83] is used for forward and backward nonunit subsumption. Prover9 is available from

<http://www.cs.unm.edu/~mccune/prover9/>

Expected Competition Performance

Some of the strategy development for CASC was done by experimentation with the CASC-2004 competition “selected” problems. (Prover9 has not yet been run on other TPTP problems.) Prover9 is unlikely to challenge the CASC leaders, because (1) extensive testing and tuning over TPTP problems has not been done, (2) theories (e.g., ring, combinatory logic, set theory) are not recognized, (3) term orderings and symbol precedences are not fine-tuned, and (4) multiple searches with differing strategies are not run.

Finishes in the middle of the pack are anticipated in all categories in which Prover9 competes.

7.21 PS-E 1.0

Daniel Kuehlwein¹, Josef Urban¹, Stephan Schulz²

¹Radboud University Nijmegen, The Netherlands, ²Technische Universität München, Germany

Architecture

PS-E (Premise Selector - E) [53] is a learning based metasystem for large theories where symbol and formula names are used consistently. It uses the MOR [53] algorithm for learning and is based upon E. PS-E learns how symbols and terms relate to used premises. Given a new conjecture, PS-E extracts its symbols, terms and subterms and predicts which premises are needed to solve it. When PS-E finds new proofs, it updates its learning algorithm to produce even more accurate predictions.

Strategies

The basic strategy is to run ATPs on problems, then use the machine learner to learn axiom relevance for conjectures from solutions, and use the most relevant axioms for next ATP attempts. This is iterated, using different timelimits and axiom limits. The learning parameters (regularization, the gaussian kernel parameters sigma, thresholds) are optimized via cross validation on the 1000 example problems and solutions.

Implementation

PS-E is implemented in python, using the numpy and scipy libraries. It is available from

<http://www.cs.ru.nl/~kuehlwein/>

Expected Competition Performance

Based on the performance of the MOR algorithms in [53] PS-E is expected to be very competitive.

7.22 Satallax 2.1

Chad E. Brown
Saarland University, Germany

Architecture

Satallax [32] is an automated theorem prover for higher-order logic. The particular form of higher-order logic supported by Satallax is Church's simple type theory with extensionality and choice operators. The SAT solver MiniSat [40] is responsible for much of the search for a proof.

The theoretical basis of search is a complete ground tableau calculus for higher-order logic [35] with a choice operator [7]. A problem is given in the THF format. A branch is formed from the axioms of the problem and the negation of the conjecture (if any is given). From this point on, Satallax tries to determine unsatisfiability or satisfiability of this branch.

Satallax progressively generates higher-order formulae and corresponding propositional clauses [32]. These formulae and propositional clauses correspond to instances of the tableau rules. Satallax uses the SAT solver MiniSat as an engine to test the current set of propositional clauses for unsatisfiability. If the clauses are unsatisfiable, then the original branch is unsatisfiable. If there are no quantifiers at function types, the generation of higher-order formulae and corresponding clauses may terminate [34, 34]. In such a case, if MiniSat reports the final set of clauses as satisfiable, then the original set of higher-order formulae is satisfiable (by a standard model in which all types are interpreted as finite sets).

Strategies

There are a number of flags that control the order in which formulas and instantiation terms are considered and propositional clauses are generated. Other flags activate some optional extensions to the basic proof procedure. A collection of flag settings is called a mode. Approximately 250 modes have been tried so far. Regardless of the mode, the search procedure is sound and complete for higher-order logic with choice. This implies that if search terminates with a particular mode, then we can conclude that the original set of formulae is unsatisfiable or satisfiable.

A strategy schedule is an ordered collection of modes with information about how much time the mode should be allotted. Satallax tries each of the modes for a certain amount of time sequentially. Satallax 2.1 has eight strategy schedules which were determined through experimentation using the THF problems in version 5.1.0 of the TPTP library. One of these eight strategy schedules is chosen based on the amount of time Satallax is given to solve the problem. For example, if Satallax is given 180 seconds to solve the problem, then a schedule with 38 modes is chosen.

Implementation

Satallax 2.1 is implemented in OCaml. A foreign function interface is used to interact with MiniSat. Satallax is available from

<http://satallax.com>

Expected Competition Performance

Satallax 2.1 is the CASC-23 THF division winner.

7.23 SPASS+T 2.2.14

Uwe Waldmann¹, Stephan Zimmer²

¹Max-Planck-Institut für Informatik, Germany, ²AbsInt GmbH, Germany

Architecture

SPASS+T is an extension of the superposition-based theorem prover SPASS that integrates algebraic knowledge into SPASS in three complementary ways: by passing derived formulas to an external SMT procedure (currently Yices or CVC3), by adding standard axioms, and by built-in arithmetic simplification and inference rules. A first version of the system has been described in [78]. In the current version, a much more sophisticated coupling of the SMT procedure has been added [128].

Strategies

Standard axioms and built-in arithmetic simplification and inference rules are integrated into the standard main loop of SPASS. Inferences between standard axioms are excluded, so the user-supplied formulas are taken as set of support. The external SMT procedure runs in parallel in a separate process, leading occasionally to non-deterministic behaviour.

Implementation

SPASS+T is implemented in C. The system is available from

<http://www.mpi-inf.mpg.de/~uwe/software/#TSPASS>

Expected Competition Performance

SPASS+T 2.2.14 is the CASC-23 TFA division winner.

7.24 SPASS+T 2.2.16

Uwe Waldmann¹, Stephan Zimmer²

¹Max-Planck-Institut für Informatik, Germany, ²AbsInt GmbH, Germany

Architecture

SPASS+T is an extension of the superposition-based theorem prover SPASS that integrates algebraic knowledge into SPASS in three complementary ways: by passing derived formulas to an external SMT procedure (currently Yices or CVC3), by adding standard axioms, and by built-in arithmetic simplification and inference rules. A first version of the system has been described in [78]. ; later a much more sophisticated coupling of the SMT procedure has been added [128].

Strategies

Standard axioms and built-in arithmetic simplification and inference rules are integrated into the standard main loop of SPASS. Inferences between standard axioms are excluded, so the user-supplied formulas are taken as set of support. The external SMT procedure runs in parallel in a separate process, leading occasionally to non-deterministic behaviour.

Implementation

SPASS+T is implemented in C. The system is available from

<http://www.mpi-inf.mpg.de/~uwe/software/#TSPASS>

Expected Competition Performance

SPASS+T 2.2.14 has been the winner of the TFA division of last CASC; we expect a similar performance in CASC 2012.

7.25 STP 1.0

Adam Pease¹, Stephan Schulz²

¹Articulate Software, USA, ²Technische Universität München, Germany

Architecture

STP is a minimalistic resolution theorem prover designed for educational purposes. The goal is to present a simple and clean series of provers starting with the most basic functional prover that can process TPTP problems, and gradually extending the system to implement more and more state-of-the-art techniques. STP is open source and it is hoped that it will provide a basis for many new derivative provers to participate in CASC.

The current system implements a complete calculus based on binary resolution.

Strategies

It includes subsumption and tautology elimination as simplification techniques, and a few simple strategies for axiom selection.

Implementation

It supports both CNF translation and proof object generation and output in TSTP format. STP is available in both Python and Java implementations.

Expected Competition Performance

Performance is expected to be poor compared to modern systems.

7.26 SuperZenon 0.0.1

David Delahaye¹, Mélanie Jacquél²

¹CEDRIC/CNAM, France, ²Siemens IC-MOL, France

Architecture

SuperZenon is an experimental extension of the Zenon [28] automated theorem prover, using the principles of superdeduction, among which the theory is used to enrich the deduction system with new deduction rules.

Superdeduction [29] is a variant of deduction modulo [39], which is an extension of logical deduction systems consisting in canonically adding ad hoc deduction rules translating axiomatic theories. This has several advantages:

- Proofs are shorter, more readable, and close to the mathematical reasoning.

- Automated proof search may speed up in such systems as some systematic parts of the proofs are "compiled" into superdeduction rules.

A version of SuperZenon has been instantiated for the set theory of the B method [1]. This allows us to provide another prover to Atelier B, which can be used to verify B proof rules in particular. This version of SuperZenon has been successfully applied (with significant speed-ups both in terms of proof time and proof size) to the verification of B proof rules coming from the database maintained by Siemens IC-MOL [49].

Strategies

The strategy of SuperZenon relies on the automatic orientation of some axioms of the theory. The axioms that can be oriented are either equivalences or implications where at least one member is atomic. The axioms which cannot be oriented are left as axioms. It should be noted that the preservation of completeness after this transformation of a part of the theory cannot be ensured in general.

Implementation

The implementation of SuperZenon relies on the one of Zenon, and has been realized thanks to the ability of Zenon to extend its core of deductive rules to match specific requirements. The compilation of superdeduction rules is performed using the proof search method of Zenon with a subset of rules. The implementation of Super Zenon is available from

<http://cedric.cnam.fr/~delahaye/super-zenon/>

Expected Competition Performance

SuperZenon is expected to improve the performance of Zenon in each division in which Zenon competes (i.e., FOFT and FOF).

7.27 TPS 3.120601S1b

Chad E. Brown¹, Peter Andrews²

¹Saarland University, Germany, ²Carnegie Mellon University, USA

Architecture

TPS is a higher-order theorem proving system that has been developed over several decades under the supervision of Peter B. Andrews with substantial work by Eve Longini Cohen, Dale A. Miller, Frank Pfenning, Sunil Issar, Carl Klapper, Dan Nesmith, Hongwei Xi, Matthew Bishop, Chad E. Brown, Mark Kaminski, Rémy Chrétien and Cris Perdue. TPS can be used to prove theorems of Church's type theory automatically, interactively, or semi-automatically [4, 5]. When searching for a proof, TPS first searches for an expansion proof [59] or an extensional expansion proof [31] of the theorem. Part of this process involves searching for acceptable matings [2]. Using higher-order unification, a pair of occurrences of subformulae (which are usually literals) is mated appropriately on each vertical path through an expanded form of the theorem to be proved. The expansion proof thus obtained is then translated [77] without further search into a proof of the theorem in natural deduction style.

Strategies

Strategies used by TPS in the search process include:

- Re-ordering conjunctions and disjunctions to alter the way paths through the formula are enumerated.
- The use of primitive substitutions and gensubs [3].
- Path-focused duplication [48].
- Dual instantiation of definitions, and generating substitutions for higher-order variables which contain abbreviations already present in the theorem to be proved [23].
- Component search [22].
- Generating and solving set constraints [30].
- Generating connections using extensional and equational reasoning [31].

Implementation

TPS has been developed as a research tool for developing, investigating, and refining a variety of methods of searching for expansion proofs, and variations of these methods. Its behavior is controlled by hundreds of flags. A set of flags, with values for them, is called a mode. The strategy of the current version of TPS consists of 71 modes. When searching for a proof in automatic mode, TPS tries each of these modes in turn for a specified amount of time (at least 1 second and at most 41 seconds). If TPS succeeds in finding an expansion proof, it translates the expansion proof to a natural deduction proof. This final step ensures that TPS will not incorrectly report that a formula has been proven. TPS is implemented in Common Lisp, and is available from

<http://gtps.math.cmu.edu/tps.html>

Expected Competition Performance

TPS was the CASC-22 THF division winner in 2009. In the two CASC competitions since 2009 TPS has come in last behind Isabelle, LEO-II and Satallax. There have been no changes to the code of TPS since last year. However, TPS has been run with many modes on the THF problems from the TPTP. This information has been used to generate a new strategy schedule which should make TPS more competitive.

7.28 Vampire-LTB 1.8

Krystof Hoder, Andrei Voronkov
University of Manchester, United Kingdom

Architecture

Vampire 1.8, is an automatic theorem prover for first-order classical logic. It consists of a shell and a kernel. The kernel implements the calculi of ordered binary resolution and superposition for handling equality. The splitting rule in kernel adds propositional parts to clauses, which are manipulated using binary decision diagrams and a SAT solver. A number of standard redundancy criteria and simplification techniques are used for pruning the search space: subsumption, tautology deletion, subsumption resolution and rewriting by ordered unit equalities. The reduction ordering is the Knuth-Bendix Ordering.

Substitution tree and code tree indexes are used to implement all major operations on sets of terms, literals and clauses. Although the kernel of the system works only with clausal normal form, the shell accepts a problem in the full first-order logic syntax, clausifies it and performs

a number of useful transformations before passing the result to the kernel. Also the axiom selection algorithm Sine [47] can be enabled as part of the preprocessing.

When a theorem is proved, the system produces a verifiable proof, which validates both the clausification phase and the refutation of the CNF.

Strategies

The Vampire 1.8 kernel provides a fairly large number of options for strategy selection. The most important ones are:

1. Choice of the main procedure:
 - Limited Resource Strategy
 - DISCOUNT loop
 - Otter loop
 - Goal oriented mode based on tabulation
2. A variety of optional simplifications.
3. Parameterized reduction orderings.
4. A number of built-in literal selection functions and different modes of comparing literals.
5. Age-weight ratio that specifies how strongly lighter clauses are preferred for inference selection.
6. Set-of-support strategy.

Implementation

Vampire 1.8 is implemented in C++.

Expected Competition Performance

Vampire-LTB 1.8 is the CASC-23 LTB division winner.

7.29 Vampire 2.6

Krystof Hoder, Andrei Voronkov
University of Manchester, England

Architecture

Vampire 2.6 is an automatic theorem prover for first-order classical logic. It consists of a shell and a kernel. The kernel implements the calculi of ordered binary resolution and superposition for handling equality. It also implements the Inst-gen calculus. The splitting rule in kernel adds propositional parts to clauses, which are being manipulated using binary decision diagrams and a SAT solver. A number of standard redundancy criteria and simplification techniques are used for pruning the search space: subsumption, tautology deletion, subsumption resolution and rewriting by ordered unit equalities. The reduction ordering is the Knuth-Bendix Ordering.

Substitution tree and code tree indexes are used to implement all major operations on sets of terms, literals and clauses. Although the kernel of the system works only with clausal normal form, the shell accepts a problem in the full first-order logic syntax, clausifies it and performs a number of useful transformations before passing the result to the kernel. Also the axiom selection algorithm Sine [47] can be enabled as part of the preprocessing.

When a theorem is proved, the system produces a verifiable proof, which validates both the clausification phase and the refutation of the CNF.

Strategies

The Vampire 2.6 kernel provides a fairly large number of options for strategy selection. The most important ones are:

- Choice of the main procedure:
 - Limited Resource Strategy
 - DISCOUNT loop
 - Otter loop
 - Goal oriented mode based on tabulation
 - Instantiation using the Inst-gen calculus
- A variety of optional simplifications.
- Parameterized reduction orderings.
- A number of built-in literal selection functions and different modes of comparing literals.
- Age-weight ratio that specifies how strongly lighter clauses are preferred for inference selection.
- Set-of-support strategy.

Implementation

Vampire 2.6 is implemented in C++.

Expected Competition Performance

We expect Vampire 2.6 to outperform the last year's Vampire 1.8, especially in the satisfiability-checking divisions.

7.30 Zenon 0.7.1

Damien Doligez
INRIA, France

Architecture

Zenon 0.7.1 [28] is a theorem prover based on a proof-confluent version of analytic tableaux. It uses all the usual tableau rules for first-order logic with a rule-based handling of equality.

Zenon outputs formal proofs that can be checked by Coq or Isabelle.

Strategies

Zenon is a fully automatic black-box design with no user-selectable strategies.

Implementation

Zenon is implemented in OCaml. Its most interesting data structure is the representation of first-order formulas and terms: they are hash-consed modulo alpha conversion.

Zenon is available from

<http://zenon-prover.org>

Expected Competition Performance

The first-order reasoning part of Zenon has not changed much since 2010, and the handling of equality is still really bad, so Zenon is again expected to rank in the lower half of the field.

The FOF division now includes some CNF problems; this will probably hinder Zenon's performance somewhat.

8 Conclusion

The CADE-23 ATP System Competition was the sixteenth large scale competition for classical logic ATP systems. The organizer believes that CASC fulfills its main motivations: stimulation of research, motivation for improving implementations, evaluation of relative capabilities of ATP systems, and providing an exciting event. Through the continuity of the event and consistency in the the reporting of the results, performance comparisons with previous and future years are easily possible. The competition provides exposure for system builders both within and outside of the community, and provides an overview of the implementation state of running, fully automatic, classical logic, ATP systems.

References

- [1] J-R. Abrial. *The B-Book, Assigning Programs to Meanings*. Cambridge University Press, 1996.
- [2] P.B. Andrews. Theorem Proving via General Matings. *Journal of the ACM*, 28(2):193–214, 1981.
- [3] P.B. Andrews. On Connections and Higher-Order Logic. *Journal of Automated Reasoning*, 5(3):257–291, 1989.
- [4] P.B. Andrews, M. Bishop, S. Issar, Nesmith. D., F. Pfenning, and H. Xi. TPS: A Theorem-Proving System for Classical Type Theory. *Journal of Automated Reasoning*, 16(3):321–353, 1996.
- [5] P.B. Andrews and C.E. Brown. TPS: A Hybrid Automatic-Interactive System for Developing Proofs. *Journal of Applied Logic*, 4(4):367–395, 2006.
- [6] L. Bachmair and H. Ganzinger. Equational Reasoning in Saturation-Based Theorem Proving. In W. Bibel and P.H. Schmitt, editors, *Automated Deduction, A Basis for Applications*, volume I Foundations - Calculi and Methods of *Applied Logic Series*, pages 352–397. Kluwer Academic Publishers, 1998.
- [7] J. Backes and C.E. Brown. Analytic Tableaux for Higher-Order Logic with Choice. *Journal of Automated Reasoning*, 47(4):451–479, 2011.
- [8] C. Barrett and C. Tinelli. CVC3. In W. Damm and H. Hermanns, editors, *Proceedings of the 19th International Conference on Computer Aided Verification*, number 4590 in Lecture Notes in Computer Science, pages 298–302. Springer-Verlag, 2007.
- [9] R. Barrett, C. Nieuwenhuis, A. Oliveras, and C. Tinelli. Splitting on Demand in SAT Modulo Theories. In M. Hermann and A. Voronkov, editors, *Proceedings of the 13th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning*, number 4246 in Lecture Notes in Artificial Intelligence, pages 512–526, 2006.
- [10] P. Baumgartner, A. Fuchs, and C. Tinelli. Darwin - A Theorem Prover for the Model Evolution Calculus. In G. Sutcliffe, S. Schulz, and T. Tammet, editors, *Proceedings of the Workshop on Empirically Successful First Order Reasoning, 2nd International Joint Conference on Automated Reasoning*, 2004.
- [11] P. Baumgartner, U. Furbach, and I. Niemelä. Hyper Tableaux. In J. Alferes, L. Pereira, and E. Orłowska, editors, *Proceedings of JELIA'96: European Workshop on Logic in Artificial Intelligence*, number 1126 in Lecture Notes in Artificial Intelligence, pages 1–17. Springer-Verlag, 1996.
- [12] P. Baumgartner, U. Furbach, and B. Pelzer. Hyper Tableaux with Equality. In F. Pfenning, editor, *Proceedings of the 21st International Conference on Automated Deduction*, number 4603 in Lecture Notes in Artificial Intelligence, pages 492–507. Springer-Verlag, 2007.
- [13] P. Baumgartner, B. Pelzer, and C. Tinelli. Model Evolution with Equality - Revised and Implemented. *Journal of Symbolic Computation*, page To appear, 2011.

- [14] P. Baumgartner and C. Tinelli. The Model Evolution Calculus. In F. Baader, editor, *Proceedings of the 19th International Conference on Automated Deduction*, number 2741 in Lecture Notes in Artificial Intelligence, pages 350–364. Springer-Verlag, 2003.
- [15] P. Baumgartner and C. Tinelli. The Model Evolution Calculus with Equality. In R. Nieuwenhuis, editor, *Proceedings of the 20th International Conference on Automated Deduction*, number 3632 in Lecture Notes in Artificial Intelligence, pages 392–408. Springer-Verlag, 2005.
- [16] C. Benzmüller. Extensional Higher-order Paramodulation and RUE-Resolution. In H. Ganzinger, editor, *Proceedings of the 16th International Conference on Automated Deduction*, number 1632 in Lecture Notes in Artificial Intelligence, pages 399–413. Springer-Verlag, 1999.
- [17] C. Benzmüller and M. Kohlhase. LEO - A Higher-Order Theorem Prover. In C. Kirchner and H. Kirchner, editors, *Proceedings of the 15th International Conference on Automated Deduction*, number 1421 in Lecture Notes in Artificial Intelligence, pages 139–143. Springer-Verlag, 1998.
- [18] C. Benzmüller, L. Paulson, F. Theiss, and A. Fietzke. LEO-II - A Cooperative Automatic Theorem Prover for Higher-Order Logic. In P. Baumgartner, A. Armando, and D. Gilles, editors, *Proceedings of the 4th International Joint Conference on Automated Reasoning*, number 5195 in Lecture Notes in Artificial Intelligence, pages 162–170. Springer-Verlag, 2008.
- [19] C. Benzmüller, F. Rabe, and G. Sutcliffe. THF0 - The Core TPTP Language for Classical Higher-Order Logic. In P. Baumgartner, A. Armando, and D. Gilles, editors, *Proceedings of the 4th International Joint Conference on Automated Reasoning*, number 5195 in Lecture Notes in Artificial Intelligence, pages 491–506. Springer-Verlag, 2008.
- [20] C. Benzmüller, V. Sorge, M. Jamnik, and M. Kerber. Combined Reasoning by Automated Cooperation. *Journal of Applied Logic*, 6(3):318–342, 2008.
- [21] W. Bibel. *Automated Theorem Proving*. Vieweg and Sohn, 1987.
- [22] M. Bishop. A Breadth-First Strategy for Mating Search. In H. Ganzinger, editor, *Proceedings of the 16th International Conference on Automated Deduction*, number 1632 in Lecture Notes in Artificial Intelligence, pages 359–373. Springer-Verlag, 1999.
- [23] M. Bishop and P.B. Andrews. Selectively Instantiating Definitions. In C. Kirchner and H. Kirchner, editors, *Proceedings of the 15th International Conference on Automated Deduction*, number 1421 in Lecture Notes in Artificial Intelligence, pages 365–380. Springer-Verlag, 1998.
- [24] J. Blanchette, S. Boehme, and L. Paulson. Extending Sledgehammer with SMT Solvers. In N. Bjorner and V. Sofronie-Stokkermans, editors, *Proceedings of the 23rd International Conference on Automated Deduction*, number 6803 in Lecture Notes in Artificial Intelligence, pages 116–130. Springer-Verlag, 2011.
- [25] J. Blanchette, A. Popescu, D. Wand, and C. Weidenbach. More SPASS with Isabelle. In L. Beringer and A. Felty, editors, *Proceedings of Interactive Theorem Proving 2012*, Lecture Notes in Artificial Intelligence, 2012.
- [26] W.W. Bledsoe. Splitting and Reduction Heuristics in Automatic Theorem Proving. *Artificial Intelligence*, 2:55–77, 1971.
- [27] S. Böhme and T. Nipkow. Sledgehammer: Judgement Day. In J. Giesl and R. Haehnle, editors, *Proceedings of the 5th International Joint Conference on Automated Reasoning*, number 6173 in Lecture Notes in Artificial Intelligence, pages 107–121, 2010.
- [28] R. Bonichon, D. Delahaye, and D. Doligez. Zenon : An Extensible Automated Theorem Prover Producing Checkable Proofs. In N. Dershowitz and A. Voronkov, editors, *Proceedings of the 14th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning*, number 4790 in Lecture Notes in Artificial Intelligence, pages 151–165, 2007.
- [29] P. Brauner, C. Houtmann, and C. Kirchner. Principles of Superdeduction. In L. Ong, editor, *Proceedings of the 22nd Annual IEEE Symposium on Logic in Computer Science*, pages 41–50. IEEE Computer Society Press, 2007.
- [30] C.E. Brown. Solving for Set Variables in Higher-Order Theorem Proving. In A. Voronkov, editor, *Proceedings of the 18th International Conference on Automated Deduction*, number 2392

- in Lecture Notes in Artificial Intelligence, pages 408–422. Springer-Verlag, 2002.
- [31] C.E. Brown. *Automated Reasoning in Higher-Order Logic: Set Comprehension and Extensionality in Church's Type Theory*. Number 10 in Studies in Logic: Logic and Cognitive Systems. College Publications, 2007.
 - [32] C.E. Brown. Reducing Higher-Order Theorem Proving to a Sequence of SAT Problems. In N. Bjorner and V. Sofronie-Stokkermans, editors, *Proceedings of the 23rd International Conference on Automated Deduction*, number 6803 in Lecture Notes in Artificial Intelligence, pages 147–161. Springer-Verlag, 2011.
 - [33] C.E. Brown. Satallax: An Automated Higher-Order Prover (System Description). In B. Gramlich, D. Miller, and U. Sattler, editors, *Proceedings of the 6th International Joint Conference on Automated Reasoning*, Lecture Notes in Artificial Intelligence, page To appear, 2012.
 - [34] C.E. Brown and G. Smolka. Terminating Tableaux for the Basic Fragment of Simple Type Theory. In M. Giese and A. Waaler, editors, *Proceedings of the 18th International Conference on Automated Reasoning with Analytic Tableaux and Related Methods*, number 5697 in Lecture Notes in Artificial Intelligence, pages 138–151. Springer-Verlag, 2009.
 - [35] C.E. Brown and G. Smolka. Analytic Tableaux for Simple Type Theory and its First-Order Fragment. *Logical Methods in Computer Science*, 6(2), 2010.
 - [36] K. Claessen, A. Lilliestrom, and N. Smallbone. Sort It Out with Monotonicity - Translating between Many-Sorted and Unsorted First-Order Logic. In N. Bjorner and V. Sofronie-Stokkermans, editors, *Proceedings of the 23rd International Conference on Automated Deduction*, number 6803 in Lecture Notes in Artificial Intelligence, pages 207–221. Springer-Verlag, 2011.
 - [37] K. Claessen and N. Sörensson. New Techniques that Improve MACE-style Finite Model Finding. In P. Baumgartner and C. Fermueller, editors, *Proceedings of the CADE-19 Workshop: Model Computation - Principles, Algorithms, Applications*, 2003.
 - [38] L. de Moura and N. Bjorner. Z3: An Efficient SMT Solver. In C. Ramakrishnan and J. Rehof, editors, *Proceedings of the 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, number 4963 in Lecture Notes in Artificial Intelligence, pages 337–340. Springer-Verlag, 2008.
 - [39] G. Dowek, T. Hardin, and C. Kirchner. Theorem Proving Modulo. *Journal of Automated Reasoning*, 31(1):33–72, 2003.
 - [40] N. Eén and N. Sörensson. An Extensible SAT-solver. In E. Giunchiglia and A. Tacchella, editors, *Proceedings of the 6th International Conference on Theory and Applications of Satisfiability Testing*, number 2919 in Lecture Notes in Computer Science, pages 502–518. Springer-Verlag, 2004.
 - [41] H. Ganzinger and K. Korovin. New Directions in Instantiation-Based Theorem Proving. In P. Kolaitis, editor, *Proceedings of the 18th IEEE Symposium on Logic in Computer Science*, pages 55–64. IEEE Press, 2003.
 - [42] H. Ganzinger and K. Korovin. Integrating Equational Reasoning into Instantiation-Based Theorem Proving. In J. Marcinkowski and A. Tarlecki, editors, *Proceedings of the 18th International Workshop on Computer Science Logic, 13th Annual Conference of the EACSL*, number 3210 in Lecture Notes in Computer Science, pages 71–84. Springer-Verlag, 2004.
 - [43] M. Gebser, O. Sabuncu, and T. Schaub. An Incremental Answer Set Programming Based System for Finite Model Computation. *AI Communications*, 24(2):195–212, 2011.
 - [44] M. Greiner and M. Schramm. A Probabilistic Stopping Criterion for the Evaluation of Benchmarks. Technical Report I9638, Institut für Informatik, Technische Universität München, München, Germany, 1996.
 - [45] K. Hoder, L. Kovacs, and A. Voronkov. Interpolation and Symbol Elimination in Vampire. In J. Giesl and R. Haehnle, editors, *Proceedings of the 5th International Joint Conference on Automated Reasoning*, number 6173 in Lecture Notes in Artificial Intelligence, pages 188–195, 2010.

- [46] K. Hoder, L. Kovacs, and A. Voronkov. Invariant Generation in Vampire. In P. Abdulla and R. Leino, editors, *Proceedings of the 17th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, number 6605 in Lecture Notes in Computer Science, pages 60–64. Springer-Verlag, 2011.
- [47] K. Hoder and A. Voronkov. Sine Qua Non for Large Theory Reasoning. In V. Sofronie-Stokkermans and N. Bjørner, editors, *Proceedings of the 23rd International Conference on Automated Deduction*, number 6803 in Lecture Notes in Artificial Intelligence, pages 299–314. Springer-Verlag, 2011.
- [48] S. Issar. Path-Focused Duplication: A Search Procedure for General Matings. In Swartout W. Dietterich T., editor, *Proceedings of the 8th National Conference on Artificial Intelligence*, pages 221–226. American Association for Artificial Intelligence / MIT Press, 1990.
- [49] M. Jacquél, K. Berkani, D. Delahaye, and C. Dubois. Tableaux Modulo Theories using Superdeduction: An Application to the Verification of B Proof Rules with the Zenon Automated Theorem Prover. In B. Gramlich, D. Miller, and U. Sattler, editors, *Proceedings of the 6th International Joint Conference on Automated Reasoning*, Lecture Notes in Artificial Intelligence, page To appear, 2012.
- [50] K. Korovin. iProver - An Instantiation-Based Theorem Prover for First-order Logic (System Description). In P. Baumgartner, A. Armando, and D. Gilles, editors, *Proceedings of the 4th International Joint Conference on Automated Reasoning*, number 5195 in Lecture Notes in Artificial Intelligence, pages 292–298, 2008.
- [51] K. Korovin. An Invitation to Instantiation-Based Reasoning: From Theory to Practice. In A. Podelski, A. Voronkov, and R. Wilhelm, editors, *Volume in Memoriam of Harald Ganzinger*, number 5663 in Lecture Notes in Computer Science, pages 163–166. Springer-Verlag, 2009.
- [52] K. Korovin and C. Stickel. iProver-Eq - An Instantiation-Based Theorem Prover with Equality. In J. Giesl and R. Haehnle, editors, *Proceedings of the 5th International Joint Conference on Automated Reasoning*, number 6173 in Lecture Notes in Artificial Intelligence, pages 196–202, 2010.
- [53] D. Kuehlwein, T. van Laarhoven, E. Tsvitshivadze, J. Urban, and T. Heskes. Overview and Evaluation of Premise Selection Techniques for Large Theory Mathematics. In B. Gramlich, D. Miller, and U. Sattler, editors, *Proceedings of the 6th International Joint Conference on Automated Reasoning*, Lecture Notes in Artificial Intelligence, page To appear, 2012.
- [54] R. Letz and G. Stenz. System Description: DCTP - A Disconnection Calculus Theorem Prover. In R. Gore, A. Leitsch, and T. Nipkow, editors, *Proceedings of the International Joint Conference on Automated Reasoning*, number 2083 in Lecture Notes in Artificial Intelligence, pages 381–385. Springer-Verlag, 2001.
- [55] B. Loechner. Things to Know When Implementing KBO. *Journal of Automated Reasoning*, 36(4):289–310, 2006.
- [56] D.W. Loveland. *Automated Theorem Proving : A Logical Basis*. Elsevier Science, 1978.
- [57] W.W. McCune. Solution of the Robbins Problem. *Journal of Automated Reasoning*, 19(3):263–276, 1997.
- [58] W.W. McCune. Otter 3.3 Reference Manual. Technical Report ANL/MS-C-TM-263, Argonne National Laboratory, Argonne, USA, 2003.
- [59] D. Miller. A Compact Representation of Proofs. *Studia Logica*, 46(4):347–370, 1987.
- [60] R. Nieuwenhuis, A. Oliveras, and C. Tinelli. Solving SAT and SAT Modulo Theories: from an Abstract Davis-Putnam-Logemann-Loveland Procedure to DPLL(T). *Journal of the ACM*, 53(6):937–977, 2006.
- [61] T. Nipkow. Equational Reasoning in Isabelle. *Science of Computer Programming*, 12(2):123–149, 1989.
- [62] T. Nipkow, L. Paulson, and M. Wenzel. Isabelle/HOL: A Proof Assistant for Higher-Order Logic. <http://www.cl.cam.ac.uk/research/hvg/Isabelle/dist/Isabelle/doc/tutorial.pdf>.

- [63] J. Otten. leanCoP 2.0 and ileancop 1.2: High Performance Lean Theorem Proving in Classical and Intuitionistic Logic. In P. Baumgartner, A. Armando, and D. Gilles, editors, *Proceedings of the 4th International Joint Conference on Automated Reasoning*, number 5195 in Lecture Notes in Artificial Intelligence, pages 283–291, 2008.
- [64] J. Otten. Restricting Backtracking in Connection Calculi. *AI Communications*, 23(2-3):159–182, 2010.
- [65] J. Otten and W. Bibel. leanCoP: Lean Connection-Based Theorem Proving. *Journal of Symbolic Computation*, 36(1-2):139–161, 2003.
- [66] J. Otten and G. Sutcliffe. Using the TPTP Language for Representing Derivations in Tableau and Connection Calculi. In B. Konev, R. Schmidt, and S. Schulz, editors, *Proceedings of the Workshop on Practical Aspects of Automated Reasoning, 5th International Joint Conference on Automated Reasoning*, pages 90–100, 2010.
- [67] D. Pastre. Automatic Theorem Proving in Set Theory. *Artificial Intelligence*, 10:1–27, 1978.
- [68] D. Pastre. Muscadet : An Automatic Theorem Proving System using Knowledge and Meta-knowledge in Mathematics. *Artificial Intelligence*, 38:257–318, 1989.
- [69] D. Pastre. Automated Theorem Proving in Mathematics. *Annals of Mathematics and Artificial Intelligence*, 8:425–447, 1993.
- [70] D. Pastre. Muscadet version 2.3 : User’s Manual. <http://www.math-info.univ-paris5.fr/pastre/muscadet/manual-en.ps>, 2001.
- [71] D. Pastre. Strong and Weak Points of the Muscadet Theorem Prover. *AI Communications*, 15(2-3):147–160, 2002.
- [72] D. Pastre. Complementarity of a Natural Deduction Knowledge-based Prover and Resolution-based Provers in Automated Theorem Proving. <http://www.math-info.univ-paris5.fr/pastre/compl-NDKB-RB.pdf>, 2007.
- [73] L. Paulson. A Generic Tableau Prover and its Integration with Isabelle. *Artificial Intelligence*, 5(3):73–87, 1999.
- [74] L. Paulson and J. Blanchette. Three Years of Experience with Sledgehammer, a Practical Link between Automatic and Interactive Theorem Provers. In G. Sutcliffe, E. Ternovska, and S. Schulz, editors, *Proceedings of the 8th International Workshop on the Implementation of Logics*, page To appear, 2010.
- [75] L.C. Paulson and T. Nipkow. *Isabelle: A Generic Theorem Prover*. Number 828 in Lecture Notes in Computer Science. Springer-Verlag, 1994.
- [76] B. Pelzer and C. Wernhard. System Description: E-KRHyper. In F. Pfenning, editor, *Proceedings of the 21st International Conference on Automated Deduction*, number 4603 in Lecture Notes in Artificial Intelligence, pages 508–513. Springer-Verlag, 2007.
- [77] F. Pfenning. *Proof Transformations in Higher-Order Logic*. PhD thesis, Carnegie-Mellon University, Pittsburg, USA, 1987.
- [78] V. Prevosto and U. Waldmann. SPASS+T. In G. Sutcliffe, R. Schmidt, and S. Schulz, editors, *Proceedings of the FLoC’06 Workshop on Empirically Successful Computerized Reasoning, 3rd International Joint Conference on Automated Reasoning*, number 192 in CEUR Workshop Proceedings, pages 19–33, 2006.
- [79] A. Riazanov and A. Voronkov. The Design and Implementation of Vampire. *AI Communications*, 15(2-3):91–110, 2002.
- [80] P. Rümmer. A Constraint Sequent Calculus for First-Order Logic with Linear Integer Arithmetic. In I. Cervesato, H. Veith, and A. Voronkov, editors, *Proceedings of the 15th International Conference on Logic for Programming Artificial Intelligence and Reasoning*, number 5330 in Lecture Notes in Artificial Intelligence, pages 274–289. Springer-Verlag, 2008.
- [81] P. Rümmer. A Constraint Sequent Calculus for First-Order Logic with Linear Integer Arithmetic. In N. Bjorner and A. Voronkov, editors, *Proceedings of the 18th International Conference on Logic*

- for Programming Artificial Intelligence and Reasoning*, number 7180 in Lecture Notes in Artificial Intelligence, pages 274–289. Springer-Verlag, 2012.
- [82] S. Schulz. A Comparison of Different Techniques for Grounding Near-Propositional CNF Formulae. In S. Haller and G. Simmons, editors, *Proceedings of the 15th International FLAIRS Conference*, pages 72–76. AAAI Press, 2002.
 - [83] S. Schulz. System Abstract: E 0.81. In M. Rusinowitch and D. Basin, editors, *Proceedings of the 2nd International Joint Conference on Automated Reasoning*, number 3097 in Lecture Notes in Artificial Intelligence, pages 223–228, 2004.
 - [84] S. Schulz. Fingerprint Indexing for Paramodulation and Rewriting. In B. Gramlich, D. Miller, and U. Sattler, editors, *Proceedings of the 6th International Joint Conference on Automated Reasoning*, Lecture Notes in Artificial Intelligence, page To appear, 2012.
 - [85] N. Sultana and C. Benzmüller. Automated Reasoning in Higher-Order Logic using the TPTP THF Infrastructure. In E. Ternovska, K. Korovin, and S. Schulz, editors, *Proceedings of the 9th International Workshop on the Implementation of Logics*, 2012.
 - [86] G. Sutcliffe. Proceedings of the CADE-16 ATP System Competition. Trento, Italy, 1999.
 - [87] G. Sutcliffe. Proceedings of the CADE-17 ATP System Competition. Pittsburgh, USA, 2000.
 - [88] G. Sutcliffe. The CADE-16 ATP System Competition. *Journal of Automated Reasoning*, 24(3):371–396, 2000.
 - [89] G. Sutcliffe. Proceedings of the IJCAR ATP System Competition. Siena, Italy, 2001.
 - [90] G. Sutcliffe. The CADE-17 ATP System Competition. *Journal of Automated Reasoning*, 27(3):227–250, 2001.
 - [91] G. Sutcliffe. Proceedings of the CADE-18 ATP System Competition. Copenhagen, Denmark, 2002.
 - [92] G. Sutcliffe. Proceedings of the CADE-19 ATP System Competition. Miami, USA, 2003.
 - [93] G. Sutcliffe. Proceedings of the 2nd IJCAR ATP System Competition. Cork, Ireland, 2004.
 - [94] G. Sutcliffe. Proceedings of the CADE-20 ATP System Competition. Tallinn, Estonia, 2005.
 - [95] G. Sutcliffe. The IJCAR-2004 Automated Theorem Proving Competition. *AI Communications*, 18(1):33–40, 2005.
 - [96] G. Sutcliffe. Proceedings of the 3rd IJCAR ATP System Competition. Seattle, USA, 2006.
 - [97] G. Sutcliffe. The CADE-20 Automated Theorem Proving Competition. *AI Communications*, 19(2):173–181, 2006.
 - [98] G. Sutcliffe. Proceedings of the CADE-21 ATP System Competition. Bremen, Germany, 2007.
 - [99] G. Sutcliffe. The 3rd IJCAR Automated Theorem Proving Competition. *AI Communications*, 20(2):117–126, 2007.
 - [100] G. Sutcliffe. Proceedings of the 4th IJCAR ATP System Competition. Sydney, Australia, 2008.
 - [101] G. Sutcliffe. The CADE-21 Automated Theorem Proving System Competition. *AI Communications*, 21(1):71–82, 2008.
 - [102] G. Sutcliffe. The SZS Ontologies for Automated Reasoning Software. In G. Sutcliffe, P. Rudnicki, R. Schmidt, B. Konev, and S. Schulz, editors, *Proceedings of the LPAR Workshops: Knowledge Exchange: Automated Provers and Proof Assistants, and The 7th International Workshop on the Implementation of Logics*, number 418 in CEUR Workshop Proceedings, pages 38–49, 2008.
 - [103] G. Sutcliffe. Proceedings of the CADE-22 ATP System Competition. Montreal, Canada, 2009.
 - [104] G. Sutcliffe. The 4th IJCAR Automated Theorem Proving System Competition - CASC-J4. *AI Communications*, 22(1):59–72, 2009.
 - [105] G. Sutcliffe. The TPTP Problem Library and Associated Infrastructure. The FOF and CNF Parts, v3.5.0. *Journal of Automated Reasoning*, 43(4):337–362, 2009.
 - [106] G. Sutcliffe. Proceedings of the 5th IJCAR ATP System Competition. Edinburgh, United Kingdom, 2010.

- [107] G. Sutcliffe. The CADE-22 Automated Theorem Proving System Competition - CASC-22. *AI Communications*, 23(1):47–60, 2010.
- [108] G. Sutcliffe. Proceedings of the CADE-23 ATP System Competition. Wroclaw, Poland, 2011.
- [109] G. Sutcliffe. The 5th IJCAR Automated Theorem Proving System Competition - CASC-J5. *AI Communications*, 24(1):75–89, 2011.
- [110] G. Sutcliffe. The CADE-23 Automated Theorem Proving System Competition - CASC-23. *AI Communications*, 25(1):49–63, 2012.
- [111] G. Sutcliffe and C. Benzmüller. Automated Reasoning in Higher-Order Logic using the TPTP THF Infrastructure. *Journal of Formalized Reasoning*, 3(1):1–27, 2010.
- [112] G. Sutcliffe, S. Schulz, K. Claessen, and A. Van Gelder. Using the TPTP Language for Writing Derivations and Finite Interpretations. In U. Furbach and N. Shankar, editors, *Proceedings of the 3rd International Joint Conference on Automated Reasoning*, number 4130 in Lecture Notes in Artificial Intelligence, pages 67–81, 2006.
- [113] G. Sutcliffe and C. Suttner. The CADE-14 ATP System Competition. Technical Report 98/01, Department of Computer Science, James Cook University, Townsville, Australia, 1998.
- [114] G. Sutcliffe and C. Suttner. The CADE-18 ATP System Competition. *Journal of Automated Reasoning*, 31(1):23–32, 2003.
- [115] G. Sutcliffe and C. Suttner. The CADE-19 ATP System Competition. *AI Communications*, 17(3):103–182, 2004.
- [116] G. Sutcliffe, C. Suttner, and F.J. Pelletier. The IJCAR ATP System Competition. *Journal of Automated Reasoning*, 28(3):307–320, 2002.
- [117] G. Sutcliffe and C.B. Suttner, editors. *Special Issue: The CADE-13 ATP System Competition*, volume 18, 1997.
- [118] G. Sutcliffe and C.B. Suttner. The Procedures of the CADE-13 ATP System Competition. *Journal of Automated Reasoning*, 18(2):163–169, 1997.
- [119] G. Sutcliffe and C.B. Suttner. Proceedings of the CADE-15 ATP System Competition. Lindau, Germany, 1998.
- [120] G. Sutcliffe and C.B. Suttner. The CADE-15 ATP System Competition. *Journal of Automated Reasoning*, 23(1):1–23, 1999.
- [121] G. Sutcliffe and C.B. Suttner. Evaluating General Purpose Automated Theorem Proving Systems. *Artificial Intelligence*, 131(1-2):39–54, 2001.
- [122] C.B. Suttner and G. Sutcliffe. The CADE-14 ATP System Competition. *Journal of Automated Reasoning*, 21(1):99–134, 1998.
- [123] E. Torlak and D. Jackson. Kodkod: A Relational Model Finder. In O. Grumberg and M. Huth, editors, *Proceedings of the 13th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, number 4424 in Lecture Notes in Computer Science, pages 632–647. Springer-Verlag, 2007.
- [124] J. Ullman. *Principles of Database and Knowledge-Base Systems*. Computer Science Press, Inc., 1989.
- [125] J. Urban. MaLAREa: a Metasystem for Automated Reasoning in Large Theories. In J. Urban, G. Sutcliffe, and S. Schulz, editors, *Proceedings of the CADE-21 Workshop on Empirically Successful Automated Reasoning in Large Theories*, number 257 in CEUR Workshop Proceedings, pages 45–58, 2007.
- [126] J. Urban, G. Sutcliffe, P. Pudlak, and J. Vyskocil. MaLAREa SG1: Machine Learner for Automated Reasoning with Semantic Guidance. In P. Baumgartner, A. Armando, and D. Gilles, editors, *Proceedings of the 4th International Joint Conference on Automated Reasoning*, number 5195 in Lecture Notes in Artificial Intelligence, pages 441–456. Springer-Verlag, 2008.
- [127] C. Wernhard. System Description: KRHyper. Technical Report Fachberichte Informatik 14–2003, Universität Koblenz-Landau, Koblenz, Germany, 2003.

- [128] S. Zimmer. Intelligent Combination of a First Order Theorem Prover and SMT Procedures. Master's thesis, Saarland University, Saarbruecken, Germany, 2007.