

CASC-V9

CASC-V9

CASC-V9

CASC-V9

Proceedings of the 9th IJCAR ATP System Competition (CASC-J9)

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 other useful sources of test problems, and specified time limits on solution attempts. The 9th IJCAR ATP System Competition (CASC-J9) was held on 14th July 2018. The design of the competition and its rules, and information regarding the competing systems, are provided in this report.

1 Introduction

The CADE and IJCAR 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 and IJCAR conference. CASC-J9 was held on 14th July 2018, as part of the 9th International Joint Conference on Automated Reasoning (IJCAR 2018)¹, which in turn was part of the Federated Logic Conference 2018, in Oxford, United Kingdom. It was the twenty-third competition in the CASC series [121, 127, 124, 81, 83, 120, 118, 119, 88, 90, 92, 94, 97, 99, 101, 103, 105, 107, 109, 126, 111, 113].

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 [114] and other useful sources of test problems, and
- specified time limits on solution attempts.

Twenty-three ATP system versions, listed in Table 1 and 2, entered into the various competition and demonstration divisions. The winners of the CASC-26 (the previous CASC) divisions were automatically entered into the corresponding demonstration 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 [121, 122, 117, 123, 79, 80, 82, 84, 85, 86, 87, 89, 91, 93, 96, 98, 100, 102, 104, 106, 108, 110, 112]. Important

¹CADE was a constituent conference of IJCAR, hence CASC-“J9”.

ATP System	Divisions	Entrant (Associates)	Entrant's Affiliation
CSE 1.0	FOF	Feng Cao (Yang Xu, Jun Liu, Shuwei Chen, Xiaomei Zhong, Peng Xu, Qinghua Liu, Huimin Fu, Xiaodong Guan, Zhenming Song)	Southwest Jiaotong University
CSE 1.1	FOF	Feng Cao (Yang Xu, Jun Liu, Shuwei Chen, Xingxing He, Xiaomei Zhong, Peng Xu, Qinghua Liu, Huimin Fu, Jian Zhong, Guanfeng Wu, Xiaodong Guan, Zhenming Song)	Southwest Jiaotong University
CSE E 1.0	FOF	Feng Cao (Yang Xu, Stephan Schulz, Jun Liu, Shuwei Chen, Xingxing He, Xiaomei Zhong, Peng Xu, Qinghua Liu, Huimin Fu, Jian Zhong, Guanfeng Wu, Xiaodong Guan, Zhenming Song)	Southwest Jiaotong University
CVC4 1.6pre	TPA FOF FNT	Andrew Reynolds (Clark Barrett, Cesare Tinelli)	University of Iowa
E 2.2pre	FOF FNT EPR LTB	Stephan Schulz	DHBW Stuttgart
Geo-III 2018C	FOF FNT EPR	Hans de Nivelle	Nazarbayev University
Grackle 0.1	LTB	Jan Jakubuv (Cezary Kaliszzyk, Stephan Schulz, Josef Urban)	Czech Technical University in Prague
iProver 2.6	EPR (demo)	CASC	CASC-26 winner
iProver 2.8	FOF FNT EPR LTB	Konstantin Korovin (Julio Cesar Lopez Hernandez)	University of Manchester
iProverModulo 2.5-0.1	FOF SLH	Guillaume Burel	University Paris-Saclay
leanCoP 2.2	FOF	Jens Otten	University of Oslo
LEO-II 1.7.0	THF	Alexander Steen (Christoph Benzmtiller)	Freie Universität Berlin
Leo-III 1.3	THF FOF EPR (demo)	Alexander Steen (Christoph Benzmtiller, Max Wisniewski)	Freie Universität Berlin

Table 1: The ATP systems and entrants

ATP System	Divisions	Entrants (Associates)	Entrant's Affiliation
MaLArea 0.7	LTB	Josef Urban (Jan Jakubuv, Cezary Kaliszyk, Stephan Schulz)	Czech Technical University in Prague
nanoCoP 1.1	FOF	Jens Otten	University of Oslo
Princess 170717	TFA	Philipp Rümmer	Uppsala University
Prover9 2009-11A	FOF (demo)	CASC (William McCune, Bob Veroff)	CASC fixed point
Satallax 3.2	THF	CASC	CASC-26 winner
Satallax 3.3	THF	Michael Färber (Chad Brown)	Universität Innsbruck
Twee 2.2	FOF	Nick Smallbone (Koen Claessen)	Chalmers University of Technology
Vampire 4.0	LTB (demo)	CASC	CASC-26 winner
Vampire 4.1	TFA FNT (demo)	CASC	CASC-26 winner
Vampire 4.2	FOF (demo)	CASC	CASC-26 winner
Vampire 4.3	TFA FOF FNT EPR LTB LTB	Giles Reger (Martin Suda, Andrei Voronkov, Evgeny Kotelnikov, Martin Riener, Laura Kovacs)	University of Manchester

Table 2: The ATP systems and entrants, continued

changes for this CASC were:

- The TFR problem category of the TFA division was put into a hiatus state.
- The SLH division was not run.
- The previous year’s winners were placed into the demonstration division, to ensure that one of the new systems wins.

The competition organizer was Geoff Sutcliffe. CASC is overseen by a panel of knowledgeable researchers who are not participating in the event. The CASC-26 panel members were Martin Giese, Aart Middeldorp, and Florian Rabe. The competition was run on computers provided by StarExec at the University of Iowa. The CASC-J9 web site provides access to resources used before, during, and after the event: <http://www.tptp.org/CASC/J9>

The CASC rules, specifications, and deadlines are absolute. Only the panel has the right to make exceptions. 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 can 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 divided into 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 ranked. Some divisions are further divided into problem categories, which makes 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 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.

Trophies are awarded to the division winners. Additionally in CASC-J9, Jasmin Blanchette at the Vrije Universiteit Amsterdam contributed a travel prize for the THF division. The winner was invited to visit Jasmin’s team at the university for up to was week. The (European segments of) travel and the hotel expenses were covered by the Matryoshka project.

The **THF** division: Typed Higher-order Form theorems (axioms with a provable conjecture). 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). The TFA division has two problem categories:

- The **TFI** category: TFA with only Integer arithmetic
- The **TFE** category: TFA with only rEal arithmetic

The **FOF** division: First-Order Form theorems (axioms with a provable conjecture). The FOF division has two problem categories:

- The **FNE** category: FOF with No Equality
- The **FEQ** category: FOF with Equality

The **FNT** division: First-order form Non-Theorems (axioms with a countersatisfiable conjecture, and satisfiable axiom sets). The FNT division has two problem categories:

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

The **EPR** division: Effectively Propositional clause normal form theorems and non-theorems (clause sets). *Effectively propositional* means that the problems are syntactically non-propositional but are known to be reducible to propositional problems, e.g., CNF problems that have no functions with arity greater than zero. 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 theorems (axioms with a provable conjecture) from Large Theories, presented in Batches. A large theory has many functors and predicates, and many axioms of which typically only a few are required for the proof of a theorem. Problems in a batch all use a common core set of axioms, and the problems in a batch are given to the ATP system all at once. Each problem category is accompanied by a set of training problems and their solutions, taken from the same source as the competition problems, that can be used for tuning and training during (typically at the start of) the competition. In CASC-J9 the LTB division had one problem category, which remained a secret until the day of CASC (to ensure there was no pre-tuning).

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, the system is a previous winner, 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 trophies or prizes are awarded.

3 Infrastructure

3.1 Computers

The computers had:

- Four (a quad-core chip) Intel(R) Xeon(R) E5-2609, 2.40GHz CPUs
- 128GB memory
- The Red Hat Enterprise Linux Server release 7.2 (Maipo) operating system, kernel 3.10.0-514.16.1.el7.x86_64.

One ATP system runs on one CPU at a time. Systems can use all the cores on the CPU (which is advantageous in the divisions where a wall clock time limit is used).

3.2 Problems

3.2.1 Problem Selection

The problems for the THF, TFA, FOF, FNT, and EPR divisions were taken from the TPTP Problem Library [114], version v7.2.0. The TPTP version used for CASC is released after the competition has started, so that new problems have not been seen by the entrants. The problems have to meet certain criteria to be eligible for selection. The problems used are randomly selected from the eligible problems based on a seed supplied by the competition panel.

- The TPTP tags problems that designed specifically to be suited or ill-suited to some ATP system, calculus, or control strategy as *biased*, and they are excluded from the competition.
- The problems must be syntactically non-propositional.
- The TPTP uses system performance data in the Thousands of Solutions from Theorem Provers (TSTP) solution library to compute problem difficulty ratings in the range 0.00 (easy) to 1.00 (unsolved) [125]. 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 if there are not enough problems with ratings in that range. Systems can be submitted before the competition so that their performance data is used for computing the problem ratings.
- The selection is constrained so that no division or category contains an excessive number of very similar problems [81].
- The selection is biased to select problems that are new in the TPTP version used, until 50% of the problems in each problem category have been selected, after which random selection (from old and new problems) continues. The number of new problems used depends on how many new problems are eligible and the limitation on very similar problems.

The problems for the LTB division are taken from various sources, with each problem category being based on one source. The process for selecting problems depends on the problem source. Entrants are expected to honestly not use publicly available problem sets for tuning before the competition.

3.2.2 Number of Problems

In the TPTP-based divisions, 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 [25] (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 are used in determining the time limit imposed on each solution attempt - see Section 3.3. The numbers of problems to be used in each division of the competition are 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 the divisions, and the time limit imposed on each solution attempt, according to the following relationship:

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

It is a lower bound on the 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, after taking into account the limitation on very similar problems, determined according to the judgement of the competition organizers.

In the LTB division the number of problems in each problem category is determined by the number of problems in the corresponding problem set. In CASC-J9, the one problem category had 5000 problems.

3.2.3 Problem Preparation

The problems are in TPTP format, with `include` directives. 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 in the TPTP-based divisions 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 the LTB division the formulae are not preprocessed, thus allowing the ATP systems to take advantage of natural structure that occurs in the problems.

In the TPTP-based divisions the problems are given in increasing order of TPTP difficulty rating. In the LTB division the problems are given in the natural order of their creation for the problem sets.

3.2.4 Batch Specification Files

The problems for each problem category of the LTB division are listed in a *batch specification* file, containing containing global data lines and one or more batch specifications. The global data lines are:

- A problem category line of the form
`division.category LTB.category_mnemonic`
- The name of a `.tgz` file (relative to the directory holding the batch specification file) that contains training data in the form of problems in TPTP format and one or more solutions to each problem in TSTP format, in a line of the form
`division.category.training_data tgz_file_name`

The `.tgz` file expands in place to three directories: `Axioms`, `Problems`, and `Solutions`. `Axioms` contains all the axiom files that are used in the training and competition problems. `Problems` contains the training problems. `Solutions` contains a subdirectory for each of the `Problems`, containing TPTP format solutions to the problem.

Each batch specification consists of:

- A header line `% SZS start BatchConfiguration`
- 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`

If the batch is ordered the ATP systems may not start any attempt on a problem, including reading the problem file, before ending the attempt on the preceding problem. For CASC-J9 it was

```
execution.order unordered
```

- 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-J9 it was

```
output.required Proof.
```

- 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. For CASC-J9 it was

```
limit.time.problem.wc 0
```
- The overall wall clock time limit for the batch, in a line of the form

```
limit.time.overall.wc limit_in_seconds
```
- 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 problem file names (relative to the directory holding the batch specification file), and output file names where the output for the problem must be written. The output files must be written in the directory specified as the second argument to the `starexec_run` script (the first argument is the name of the batch specification file).
- A terminator line `% SZS end BatchProblems`

3.3 Resource Limits

In the TPTP-based division, CPU and wall clock time limits are imposed for each problem. 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 is imposed, depending on the computers' memory.

In the LTB division, wall clock time limits are imposed. For each batch there might be a wall clock time limit for each problem, provided in the configuration section at the start of each batch. The minimal wall clock time limit per problem is 15s, and the maximal wall clock time limit per problem is 90s. For each batch there is an overall wall clock time limit, provided in the configuration section at the start of each batch. The overall limit is proportional to the number of problems in the batch, e.g., the batch's per-problem time 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 the end of an attempt on a problem and the starting of 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, the wall clock time taken, and whether or not a proof or model was output.

The systems are ranked in the competition divisions, from the performance data. The THF, TFA, FOF, FNT, and LTB divisions are ranked according to the number of problems solved with an acceptable proof/model output. The EPR division is ranked according to the number of problems solved, but not necessarily accompanied by a proof or model (but systems that do output proofs/models are highlighted in the presentation of results). Ties are broken according to the average time taken over problems solved (CPU time or wall clock time, depending on the type of limit in the division).

The competition panel decides whether or not the systems' proofs and models are "acceptable". 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, e.g., 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 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 reciprocal of the number of systems that solved the problem, so that if a system is the only one to solve a problem then its SotAC for the problem is 1.00, and if all the systems solve a problem their SotAC for the problem is the inverse of the number of systems. A system's overall SotAC is its average SotAC over the problems it solves.
- 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 of multiple cores. The competition ran on quad-core computers, thus the maximal core usage was 4.0.
- The *efficiency* measure combines the number of problems solved with the time taken. It is the average of the inverses of the times taken for problems solved, 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. Efficiency is computed for both CPU time and wall clock time, to measure how efficiently the systems use one core and how efficiently systems use multiple cores, respectively.

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 divisions ranked by the numbers of proofs/models output) 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, by the registration deadline. For each system an entrant must be nominated to handle all issues (e.g., installation and 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 demonstration divisions, to provide benchmarks against which progress can be judged. Prover9 2009-11A is automatically entered into the FOF division, to provide a fixed-point against which progress can be judged.

5.1 System Descriptions

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 the system is also given 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 form part of the competition proceedings.

5.2 Sample Solutions

For systems in the divisions that require proof/model output, 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. The competition panel decides whether or not proofs and models are acceptable.

Proof/model samples are required as follows:

- THF: SET0144
- TFA: DAT013=1
- FOF and LTB: SEU140+2
- FNT: NLP042+1 and SWV017+1

An explanation must be provided for any non-obvious features.

6 System Requirements

6.1 System Properties

Entrants must ensure that their systems execute in the competition environment, and have the following properties. Entrants are advised to finalize their installation packages and check these properties well in advance of the system delivery deadline. This gives the competition organizers time to help resolve any difficulties encountered.

Execution, Soundness, and Completeness

- Systems must be fully automatic, i.e., all command line switches have to be the same for all problems in each division.
- Systems' performance must be reproducible by running the system again.
- 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 THF, TFA, FOF, EPR, and LTB divisions, and theorems are submitted to the systems in the 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. 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 problems and axiom sets is not allowed. Strategies and strategy selection based on individual problems is not allowed. If machine learning procedures are used, the learning must ensure that sufficient generalization is obtained so that there is no specialization to individual problems or their solutions.

- All techniques used must be general purpose, and expected to extend usefully to new unseen problems. The precomputation and storage of information about individual problems that might appear in the competition or their solutions is not allowed. (It's OK to store information about LTB training problems.) Strategies and strategy selection based on individual problems or their solutions are not allowed. If machine learning procedures are used to tune a system, the learning must ensure that sufficient generalization is obtained so that there is no specialization to individual problems or their solutions. The system description must explain any such tuning or training that has been done. The competition panel may disqualify any system that is deemed to be problem specific rather than general purpose.

Output

- In all divisions except LTB all solution output must be to `stdout`. In the LTB division all solution output must be to the named output file for each problem, in the directory specified as the second argument to the `starexec_run` script. If multiple attempts are made on a problem in an unordered batch, each successive output file must overwrite the previous one.
- 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). For example

```
% SZS status Started for CSR075+2.p
... (system churns away, result and solution output to file)
% SZS status GaveUp for CSR075+2.p
% SZS status Ended for CSR075+2.p</PRE>
... and later in another attempt on that problem ...<PRE>
% SZS status Started for CSR075+2.p
... (system churns away, result and solution appended to file)
% SZS status Theorem for CSR075+2.p
% SZS status Ended for CSR075+2.p
```

- For each problem, the system must output a distinguished string indicating what solution has been found or that no conclusion has been reached. Systems must use the SZS ontology and standards [95] for this. For example

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

or

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

In the LTB division this line must be the last line output before the ending notification line. The line must also be output to the output file.

- When outputting proofs/models, the start and end of the proof/model must be delimited by distinguished strings. Systems must use the SZS ontology and standards for this. For example

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

```
...
```

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

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 [16] is encouraged.

Resource Usage

- Systems must be interruptible by a `SIGXCPU` signal, so that CPU time limits can be imposed, and interruptible by a `SIGALRM` signal, so that wall clock time limits 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 a system terminates of its own accord, it may not leave any temporary or intermediate output files. If a system is terminated by a `SIGXCPU` or `SIGALRM`, it may not leave any temporary or intermediate files anywhere other than in `/tmp`.
- 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.

6.2 System Delivery

Entrants must email a StarExec installation package to the competition organizers by the system delivery deadline. The installation package must be a `.tgz` file containing only the components necessary for running the system (i.e., not including source code, etc.). The entrants must also email a `.tgz` file containing the source code and any files required for building the StarExec installation package to the competition organizers by the system delivery deadline.

For systems running on entrant supplied computers in the demonstration division, entrants must email a `.tgz` file containing the source code and any files required for building the executable system to the competition organizers by the system delivery deadline.

After the competition all competition division systems' source code is made publicly 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 is controlled by StarExec. The jobs are queued onto the computers so that each CPU is running one job at a time. All attempts at the Nth problems in all the divisions and categories are started before any attempts at the (N+1)th problems.

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 CSE 1.0

Feng Cao (Yang Xu, Jun Liu, Shuwei Chen, Xiaomei Zhong, Peng Xu, Qinghua Liu, Huimin Fu, Xiaodong Guan, Zhenming Song)
 Southwest Jiaotong University, China
 Ulster University, United Kingdom (Jun Liu)

Architecture

CSE 1.0 is an automated theorem prover for first-order logic without equality mainly based on a novel inference mechanism, called as Contradiction Separation Based Dynamic Multi-Clause Synergized Automated Deduction (S-CS) [133]. This S-CS inference rule is able to handle multiple (two or more) clauses dynamically in a synergized way in one deduction step, while binary resolution is its special case. In each step, multiple clauses are selected and separated into two parts, sub-clauses, while one part of each clause is used to form a contradiction, and the disjunction of the remaining literals forms the logical consequence of the selected clauses, called contradiction separation clause (CSC).

CSE 1.0 adopts conventional factoring, equality resolution, and variable renaming. Some pre-processing techniques, including pure literal deletion and simplification based on the distance to the goal clause, and a number of standard redundancy criteria for pruning the search space: tautology deletion, subsumption (forward and backward) are applied as well.

Internally, CSE 1.0 works only with clausal normal form. E prover [77] is adopted with thanks for clausification of full first-order logic problems during preprocessing.

Strategies

CSE 1.0 has provided a number of options for strategy selection. The most important ones are:

- Clause selection. This strategy category mainly considers the times of the clauses taken part in the inference, the weights considering clause redundancy, term weight, clause complexity, the number of literals, the number of complementary predicates, the clause being a source clause or an inferred clause, and so on.
- Literal selection. This strategy category mainly considers the times of the literal taken part in the inference, literal stability, literal complexity, the number of predicates in the literal, etc.
- Weight strategy. The weights to the clauses are calculated mainly considering the times of clause taken part in the deduction process, and clause redundancy. The weights are updated dynamically during the deduction process.
- Contradiction separation clause (CSC) strategy. This strategy category mainly considers the number of literals in the CSC, the percent of the ground literals appearing in the CSC, the number of function symbol acting on the literals in a clause, and the effectiveness evaluation of the CSC.
- Control strategy of the contradiction separation deduction process. The number of literals in the CSC, and the number of clauses involved in each contradiction are dynamically changed during the deduction process.

Implementation

CSE 1.0 is implemented mainly in C++, and JAVA is used for batch problem running implementation. Shared data structure is used for constants and shared variables storage. E prover is used for clausification of FOF problems, and then TPTP2X is applied to convert the CNF format into TPTP format.

Expected Competition Performance

CSE is new to CASC, and we expect it to have an acceptable performance.

Acknowledgement

Development of CSE 1.0 has been partially supported by the National Natural Science Foundation of China (NSFC) (Grant No.61673320) and the Fundamental Research Funds for the Central Universities in China (Grant No.2682018ZT10).

7.2 CSE 1.1

Feng Cao (Yang Xu, Jun Liu, Shuwei Chen, Xingxing He, Xiaomei Zhong, Peng Xu, Qinghua Liu, Huimin Fu, Jian Zhong, Guanfeng Wu, Xiaodong Guan, Zhenming Song)
Southwest Jiaotong University, China
Ulster University, United Kingdom (Jun Liu)

Architecture

The basic inference mechanism of CSE 1.1 is similar to CSE 1.0, i.e., it is an automated theorem prover for first-order logic without equality mainly based on a novel inference mechanism, called as Contradiction Separation Based Dynamic Multi-Clause Synergized Automated Deduction (S-CS) [133], which is able to handle multiple (two or more) clauses dynamically in a synergized way in one deduction step, while binary resolution is its special case. The difference between CSE 1.0 and CSE 1.1 is that there are two S-CS deduction mechanisms in CSE 1.1, where one is called from left to right, which refers the clauses that are not in the contradiction under construction, and another is named from right to left, which considers the clauses that are already in the contradiction under construction. In addition, it supports the repeat usage of the same clause in one deduction step. These characteristics make the S-CS deduction be able to produce more unit clauses.

CSE 1.1 adopts conventional factoring, equality resolution, and variable renaming. Some pre-processing techniques, including pure literal deletion and simplification based on the distance to the goal clause, and a number of standard redundancy criteria for pruning the search space: tautology deletion, subsumption (forward and backward) are applied as well.

Internally, CSE 1.1 works only with clausal normal form. E prover [77] is adopted with thanks for clausification of full first-order logic problems during preprocessing.

Strategies

Strategies CSE 1.1 inherited most of the clause/literal selection strategy selection, while the crucial difference comes from the multiple strategy mode and some heuristic strategies. The

multiple strategy mode allows CSE 1.1 to solve the problem by trying different combination of strategies. Besides the strategies used in CSE 1.0, e.g., clause selection, literal selection, and weight strategy, there are some different strategies:

- Deduction framework. This provides two overall options for S-CS deduction: integrity deduction mode, which takes all the clauses into consideration during deduction process, and contradiction separation clause deduction mode, which considers only a subset of clauses.
- Repeat usage of clause. This strategy provides two strategies: repeat usage of axiom and repeat usage of clause.
- Contradiction separation clause strategy. Besides the CSC strategies in CSE 1.0, CSE 1.1 allows the usage of the medium CSCs during the contradiction construction process.

Implementation

CSE 1.1 is implemented mainly in C++, and JAVA is used for batch problem running implementation. Shared data structure is used for constants and shared variables storage. In addition, special data structure is designed for property description of clause, literal and term, so that it can support the multiple strategy mode. E prover is used for clausification of FOF problems, and then TPTP2X is applied to convert the CNF format into TPTP format.

Expected Competition Performance

CSE is new to CASC, and we expect it to have an acceptable performance.

Acknowledgement

Development of CSE 1.0 has been partially supported by the National Natural Science Foundation of China (NSFC) (Grant No.61673320) and the Fundamental Research Funds for the Central Universities in China (Grant No.2682018ZT10).

7.3 CSE_E 1.0

Feng Cao (Yang Xu, Stephan Schulz, Jun Liu, Shuwei Chen, Xingxing He, Xiaomei Zhong, Peng Xu, Qinghua Liu, Huimin Fu, Jian Zhong, Guanfeng Wu, Xiaodong Guan, Zhenming Song)

Southwest Jiaotong University, China

DHBW Stuttgart, Germany (Stephan Schulz)

Ulster University, United Kingdom (Jun Liu)

Architecture

CSE_E 1.0 is an automated theorem prover for first-order logic by combining CSE 1.1 and E 2.1, where CSE is based on the Contradiction Separation Based Dynamic Multi-Clause Synergized Automated Deduction (S-CS) [133] and E is based on superposition. The combination mechanism is like this: E and CSE are applied to the given problem sequentially. If either prover solves the problem, then the proof process completes. If neither CSE nor E can solve

the problem, some inferred clauses, especially unit clauses, by CSE will be fed to E as lemmas, along with the original clauses, for further proof search.

This kind of combination is expected to take advantage of both CSE and E, and produce a better performance. Concretely, CSE is able to generate a good number of unit clauses, based on the fact that unit clauses are helpful for proof search and equality handling. On the other hand, E has a good ability on equality handling.

Strategies

The strategies of CSE part of CSE_E 1.0 take the same strategies as in CSE 1.1 standalone, e.g., clause/literal selection, strategy selection, and CSC strategy. The different built-in strategies in CSE_E 1.1 are:

- Evaluation of contradiction separation clause. The evaluation is based on the number of clauses, the distance to the goal clause, and other factors.
- Deletion of contradiction separation clause. This is realized based on weights. Three weight calculation methods considering variables, functions, terms, and the time of a clause being involved in the deduction are provided.

Implementation

CSE_E 1.0 is implemented mainly in C++, and JAVA is used for batch problem running implementation. The job dispatch between CSE and E is implemented in JAVA.

Expected Competition Performance

CSE_E is new to CASC, and it is the first step trying to incorporate CSE with other theorem provers. We expect it to have a satisfying performance.

Acknowledgement

Development of CSE 1.0 has been partially supported by the National Natural Science Foundation of China (NSFC) (Grant No.61673320) and the Fundamental Research Funds for the Central Universities in China (Grant No.2682018ZT10). Thanks should also be given to Prof. Stephan Schulz for his kind permission on using his E prover that makes CSE-E possible.

7.4 CVC4 1.6pre

Andrew Reynolds
University of Iowa, USA

Architecture

CVC4 [5] is an SMT solver based on the DPLL(T) architecture [50] that includes built-in support for many theories, including linear arithmetic, arrays, bit vectors, datatypes, finite sets and strings. It incorporates approaches for handling universally quantified formulas. For problems involving free function and predicate symbols, CVC4 primarily uses heuristic approaches

based on E-matching for theorems, and finite model finding approaches for non-theorems. For problems in pure arithmetic, CVC4 uses techniques for counterexample-guided quantifier instantiation [64].

Like other SMT solvers, CVC4 treats quantified formulas using a two-tiered approach. First, quantified formulas are replaced by fresh Boolean predicates and the ground theory solver(s) are used in conjunction with the underlying SAT solver to determine satisfiability. If the problem is unsatisfiable at the ground level, then the solver answers “unsatisfiable”. Otherwise, the quantifier instantiation module is invoked, and will either add instances of quantified formulas to the problem, answer “satisfiable”, or return unknown. Finite model finding in CVC4 targets problems containing background theories whose quantification is limited to finite and uninterpreted sorts. In finite model finding mode, CVC4 uses a ground theory of finite cardinality constraints that minimizes the number of ground equivalence classes, as described in [66]. When the problem is satisfiable at the ground level, a candidate model is constructed that contains complete interpretations for all predicate and function symbols. It then adds instances of quantified formulas that are in conflict with the candidate model, as described in [67]. If no instances are added, it reports “satisfiable”.

Strategies

For handling theorems, CVC4 primarily uses conflict-based quantifier instantiation [65, 4] and E-matching. CVC4 uses a handful of orthogonal trigger selection strategies for E-matching. For handling non-theorems, CVC4 primarily uses finite model finding techniques. Since CVC4 with finite model finding is also capable of establishing unsatisfiability, it is used as a strategy for theorems as well. For problems in pure arithmetic, CVC4 uses variations of counterexample-guided quantifier instantiation [64], which select relevant quantifier instantiations based on models for counterexamples to quantified formulas. CVC4 relies on this method both for theorems in TFA and non-theorems in TFN. At the quantifier-free level, CVC4 uses standard decision procedures for linear arithmetic and uninterpreted functions, as well as heuristic approaches for handling non-linear arithmetic [68].

Implementation

CVC4 is implemented in C++. The code is available from:

<https://github.com/CVC4>

Expected Competition Performance

The first-order theorem proving and finite model finding capabilities of CVC4 have not changed much in the past year. Its support for non-linear arithmetic has been improved. It is expected that CVC4 will perform slightly better than last year.

7.5 E 2.2pre

Stephan Schulz
DHBW Stuttgart, Germany

Architecture

E 2.1 [73, 77] is a purely equational theorem prover for many-sorted 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. However, as of E 2.1, PicoSAT [12] can be used to periodically check the (on-the-fly grounded) proof state for propositional unsatisfiability.

For LTB division, a control program uses a SInE-like analysis to extract reduced axiomatizations that are handed to several instances of E. E will probably not use on-the-fly learning this year.

Strategies

Proof search in E is primarily controlled by a literal selection strategy, a clause selection heuristic, and a simplification ordering. The prover supports a large number of pre-programmed literal selection strategies. Clause selection heuristics can be constructed on the fly by combining various parameterized 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 parameterized instances of Knuth-Bendix-Ordering (KBO) and Lexicographic Path Ordering (LPO).

For CASC-J9, E implements a strategy-scheduling automatic mode. The total CPU time available is broken into several (unequal) time slices. For each time slice, the problem is classified into one of several classes, based on a number of simple features (number of clauses, maximal symbol arity, presence of equality, presence of non-unit and non-Horn clauses,...). For each class, a schedule of strategies is greedily constructed from experimental data as follows: The first strategy assigned to a schedule is the one that solves the most problems from this class in the first time slice. Each subsequent strategy is selected based on the number of solutions on problems not already solved by a preceding strategy. About 220 different strategies have been evaluated on all untyped first-order problems from TPTP 6.4.0. About 90 of these strategies are used in the automatic mode, and about 210 are used in at least one schedule.

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 [76]. Superposition and backward rewriting use fingerprint indexing [75], 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 [45, 46]. The prover and additional information are available at

<http://www.eprover.org>

Expected Competition Performance

E 2.2pre has only minor changes compared to last years pre-releases. The major change is the integration of PicoSAT, however, very few PicoSAT strategies have been evaluated. As a result, we expect performance to be similar to last years, with maybe most improvements in the EPR division. The system is expected to perform well in most proof classes, but will at best complement top systems in the disproof classes.

7.6 Geo-III 2018C

Hans de Nivelle
Nazarbayev University, Kazakhstan

Architecture

Geo III is a theorem prover for Partial Classical Logic [20], based on reduction to Kleene Logic [21]. Currently, only Chapters 4 and 5 are implemented. Since Kleene logic generalizes 2-valued logic, Geo III can take part in CASC. Apart from being 3-valued, the main differences with earlier versions of Geo are the following:

1. The Geo family of provers uses exhaustive backtracking, in combination with learning after failure. Earlier versions (before 2016) learned only conflict formulas. Geo III learns disjunctions of arbitrary width. Experiments show that this often results in shorter proofs.
2. If Geo will be ever embedded in proof assistants, these assistants will require proofs. In order to be able to provide these at the required level of detail, Geo III contains a hierarchy of proof rules that is independent of the rest of the system, and that can be modified independently.
3. In order to be more flexible in the main algorithm, recursive backtracking has been replaced by use of a stack. By using a stack, it has become possible to implement non-chronological backtracking, remove unused assumptions, or to rearrange the order of assumptions. Also, restarts are easier to implement with a stack.
4. Matching a geometric formula into a candidate model is a critical operation in Geo. Compared to previous versions, the matching algorithm has been improved theoretically, reimplemented, and is no longer a bottle neck.

As for future plans, we want to add backward simplification to the main algorithm. This involves matching between geometric formulas, which was not possible before, because we had no usable matching algorithm. We also want to reimplement proof logging, and to implement full PCL.

Strategies

Geo uses breadth-first, exhaustive model search, combined with learning. In case of branching, branches are explored in pseudo-random order. Specially for CASC, a restart strategy was added, which ensures that proof search is always restarted after 4 minutes. This was done because Geo III has no indexing. After some time, proof search becomes so inefficient that it makes no sense to continue, so that it is better to restart.

Implementation

Geo III is written in C++-14. No features outside of the standard are used. It has been tested with g++ (version 4.8.4) and with clang. The main difference with Geo 2016C is that version 2018C uses a new matching algorithm, which on average performs 100 to 1000 times better than the previous one. Geo-III is available at:

<https://cs-sst.github.io/faculty/nivelle/implementation/index>

Expected Competition Performance

We are slowly closing the gaps in Geo. We expect Geo 2018C to be better than 2016C, but the way to the top is long.

Acknowledgement

Development of Geo 2018C was supported by the Polish National Science Center (NCN) through grant number DEC-2015/17/B/ST6/01898 (Applications for Logic with Partial Functions).

7.7 Grackle 0.1

Jan Jakubuv
Czech Technical University in Prague, Czech Republic

Architecture

Grackle is a bird species found in large numbers through much of North America. Different subspecies of the grackle family evolved a different bill length. This has the effect that different subspecies feed on different nutriment and do not compete with each other. This motivates the Grackle system. Grackle 0.1 is a generalization of BliStrTune [129, 33, 34], a system for invention of complementary E prover strategies, based on ParamILS system [31]. BliStrTune was previously extended to invent Vampire strategies [32] but this is not used here. Grackle is a next step in this direction of generalization, and it is able to develop complementary strategies of an arbitrary parametrized algorithm, not only E or Vampire. In CASC-J9, however, Grackle is used only to develop E strategies and the main difference from BliStrTune is a separate invention of SinE parameters for E prover.

Strategies

The basic strategy is to divide the given time limit for the LTB category into halves. In the first half, Grackle is launched to invent a set of well-performing complementary E strategies on the provided training examples. Provided solutions of the training examples are not used. In the second half of the time limit, the invented strategies are evaluated using E prover. Hence the output solutions are in E prover output format (TPTP). The evaluation again has two stages. Firstly with a small CPU limit (e.g. 10 seconds per strategy and problem) and then the unsolved problems are evaluated with a longer limit (e.g. 60 seconds).

Implementation

The metasystem is implemented using ATPy Python library available at:

<https://github.com/ai4reason/atpy>

The Grackle system is a part of this library. Grackle itself uses ParamILS to improve an existing strategy. Grackle requires a set of reasonably performing E prover strategies to start with. These are extracted from E's auto mode and previous Grackle/BliStrTune runs on a subset of TPTP library. The code of ATPy and Grackle is released under GPL2.

Expected Competition Performance

Grackle will compete only in the LTB category. First solutions are to be expected after the training phase, that is, after about half of the overall time limit. As the problems in the LTB category are expected to be large, Grackle would normally require several days (maybe weeks) for a successful training. This is because a separate premise selection is still missing in Grackle and only the SinE algorithm implemented natively in E prover is used. Hence the expectations are humble and Grackle can only surprise.

7.8 iProver 2.6

Konstantin Korovin
University of Manchester, United Kingdom

Architecture

iProver is an automated theorem prover based on an instantiation calculus Inst-Gen [23, 40] which is complete for first-order logic. iProver combines first-order reasoning with ground reasoning for which it uses MiniSat [22] and optionally PicoSAT [12] (only MiniSat will be used at this CASC). iProver also combines instantiation with ordered resolution; see [39, 40] 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; mismatching constraints [24, 39]; global subsumption [39]; 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. Recent changes in iProver include improved preprocessing and incremental finite model finding; support for the TFF format restricted to clauses; the AIG format for hardware verification and QBF reasoning.

In the LTB and SLH divisions, iProver combines an abstraction-refinement framework [27] with axiom selection based on the SinE algorithm [30] as implemented in Vampire [43], i.e., axiom selection is done by Vampire and proof attempts are done by iProver.

Some of iProver features are summarised below.

- proof extraction for both instantiation and resolution [42],
- model representation, using first-order definitions in term algebra [42],
- answer substitutions,
- semantic filtering,
- incremental finite model finding,
- sort inference, monotonic [17] and non-cyclic [41] sorts,
- support for the TFF format restricted to clauses,
- predicate elimination [38].

Sort 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. For the LTB, SLH and FNT divisions several strategies are run in parallel.

Implementation

iProver is implemented in OCaml and for the ground reasoning uses MiniSat [22]. iProver accepts FOF, TFF and CNF formats. Vampire [43, 28] and E prover [77] are used for proof-producing classification of FOF/TFF problems, Vampire is also used for axiom selection [HV11] in the LTB/SLH divisions. iProver is available at:

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

Expected Competition Performance

iProver 2.6 is the CASC-26 EPR division winner.

7.9 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 [49]. 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

Prover9 has available many strategies; the following statements apply to CASC.

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.

For CASC Prover9 uses KBO to order terms for demodulation and for the inference rules, with a simple rule for determining symbol precedence.

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 [48]. (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 [74] is used for forward and backward nonunit subsumption. Prover9 is available from

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

Expected Competition Performance

Prover9 is the CASC fixed point, against which progress can be judged. Each year it is expected do worse than the previous year, relative to the other systems.

7.10 leanCoP 2.2

Jens Otten
University of Oslo, Norway

Architecture

leanCoP [60, 51] is an automated theorem prover for classical first-order logic with equality. It is a very compact implementation of the connection (tableau) calculus [11, 44].

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 performed in order to achieve completeness. Additional inference rules and techniques include regularity, lemmata, and restricted backtracking [52]. leanCoP uses an optimized structure-preserving transformation into clausal form [52] and a fixed strategy schedule that is controlled by a shell script.

Implementation

leanCoP is implemented in Prolog. The source code of the core prover consists only of a few lines of code. Prolog's built-in indexing mechanism is used to quickly find connections when the extension rule is applied.

leanCoP can read formulae in leanCoP syntax and in 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 is then translated into a more comprehensive output format, e.g., into a lean (TPTP-style) connection proof or into a readable text proof.

The source code of leanCoP 2.2 is available under the GNU general public license. It can be downloaded from the leanCoP website at:

`http://www.leancop.de`

The website also contains information about ileanCoP [51] and MleanCoP [54, 55], two versions of leanCoP for first-order intuitionistic logic and first-order modal logic, respectively.

Expected Competition Performance

As the prover has not changed, the performance of leanCoP 2.2 is expected to be the same as last year.

7.11 LEO-II 1.7.0

Alexander Steen
Freie Universität Berlin, Germany

Architecture

LEO-II [8], the successor of LEO [7], 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 [6]. 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 [72]. 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 uses some 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 OCaml 4, and its problem representation language is the TPTP THF language [9]. In fact, the development of LEO-II has largely paralleled the development of the TPTP THF language and related infrastructure [115]. 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 [10] 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 is not actively being developed anymore, hence there are no expected improvements to last year's CASC results.

7.12 Leo-III 1.3

Alexander Steen
Freie Universität Berlin, Germany

Architecture

Leo-III [78], the successor of LEO-II [8], is a higher-order ATP system based on extensional higher-order paramodulation with inference restrictions using a higher-order term ordering. The calculus is augmented with dedicated extensionality rules and equational simplification routines that have their intellectual roots in first-order superposition-based theorem proving. Although Leo-III is originally designed as an agent-based reasoning system, its current version utilizes one sequential saturation algorithm only. The saturation algorithm itself is a variant of the given clause loop procedure.

Leo-III heavily relies on cooperation with external (first-order) ATPs that are called asynchronously during proof search. At the moment, first-order cooperation focuses on typed first-order (TFF) systems, where CVC4 [5] and E [73, 77] are used as default external systems. Nevertheless, cooperation is not limited to first-order systems. Further TPTP/TSTP-compliant external systems (such as higher-order ATPs or counter model generators) may be included using simple command-line arguments. If the saturation procedure loop (or one of the external provers) finds a proof, the system stops, generates the proof certificate and returns the result.

Strategies

Leo-III comes with pre-defined search strategies that can be chosen manually by the user on start-up. However, currently, Leo-III supports only very naive automatic strategy scheduling that is, by default, disabled as its effectivity seems not well-examined yet. Strategies will primarily be addressed in further upcoming versions.

Implementation

Leo-III exemplarily utilizes and instantiates the associated LeoPARD system platform [132] for higher-order (HO) deduction systems implemented in Scala (currently using Scala 2.12 and running on a JVM with Java 8). The prover makes use of LeoPARD's sophisticated data structures and implements its own reasoning logic on top. A generic parser is provided that supports all TPTP syntax dialects. It is implemented using ANTLR4 and converts its produced concrete syntax tree to an internal TPTP AST data structure which is then transformed into polymorphically typed lambda terms. As of version 1.1, Leo-III supports all common TPTP dialects (CNF, FOF, TFF, THF) as well as its polymorphic variants [13, 35].

The term data structure of Leo-III uses a polymorphically typed spine term representation augmented with explicit substitutions and De Bruijn-indices. Furthermore, terms are perfectly shared during proof search, permitting constant-time equality checks between alpha-equivalent terms.

Leo-III's saturation procedure may at any point invoke external reasoning tools. To that end, Leo-III includes an encoding module which translates (polymorphic) higher-order clauses to polymorphic and monomorphic typed first-order clauses, whichever is supported by the external system. While LEO-II relied on cooperation with untyped first-order provers, Leo-III

exploits the native type support in first-order provers (TFF logic) for removing clutter during translation and, in turn, higher effectivity of external cooperation.

Leo-III is available on GitHub:

<https://github.com/leoprover/Leo-III>

Expected Competition Performance

As of Leo-III 1.1, novel cooperation schemes with typed first-order provers were used that significantly increased the reasoning capabilities of Leo-III. In version 1.3, some minor bug fixes and parameter tweaks were conducted which should improve Leo-III, at least to some extent, compared to last year's performance. Some hard problems may be solved by Leo-III's function synthesis capabilities.

7.13 MaLAREa 0.6

Josef Urban

Czech Technical University in Prague, Czech Republic

Architecture

MaLAREa 0.6 [128, 130, 37] 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-J9 will use the E prover with the BliStr(Tune) [129, 33] large-theory strategies, possibly also Prover9, Mace and Paradox. The premise selection methods will likely also use the distance-weighted k-nearest neighbor [36] 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 at:

<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 on several previous LTB competitions.

7.14 nanoCoP—1.1

Jens Otten
University of Oslo, Norway

Architecture

nanoCoP [56, 57] is an automated theorem prover for classical first-order logic with equality. It is a very compact implementation of the non-clausal connection calculus [53].

Strategies

An additional decomposition rule is added to the clausal connection calculus and the extension rule is generalized to non-clausal formulae. Open branches are selected in a depth-first way. Iterative deepening on the proof depth is performed in order to achieve completeness. Additional inference rules and techniques include regularity, lemmata, restricted backtracking, and a fixed strategy schedule that is controlled by a shell script [59].

Implementation

nanoCoP is implemented in Prolog. The source code of the core prover consists only of a few lines of code. Prolog's built-in indexing mechanism is used to quickly find connections when the extension rule is applied.

nanoCoP can read formulae in leanCoP/nanoCoP syntax and in TPTP first-order syntax. Equality axioms are automatically added if required. The nanoCoP core prover returns a compact non-clausal connection proof.

The source code of nanoCoP 1.1 is available under the GNU general public license. It can be downloaded from the nanoCoP website at:

`http://www.leancop.de/nanocop`

The provers nanoCoP-i and nanoCoP-M are version of nanoCoP for first-order intuitionistic logic and first-order modal logic, respectively. They are based on an adapted non-clausal connection calculus for non-classical logics [58].

Expected Competition Performance

nanoCoP is expected to have a better performance than leanCoP on formulae that have a nested (non-clausal) structure.

7.15 Princess 170717

Philipp Rümmer
Uppsala University, Sweden

Architecture

Princess [70, 71] is a theorem prover for first-order logic modulo linear integer arithmetic. The prover uses 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. In addition, some built-in procedures for nonlinear integer arithmetic are available.

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

For CASC, Princess will run a fixed schedule of configurations for each problem (portfolio method). Configurations determine, among others, the mode of proof expansion (depth-first, breadth-first), selection of triggers in quantified formulae, clausification, and the handling of functions. The portfolio was chosen based on training with 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 used.

Princess is available from:

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

Expected Competition Performance

Princess should perform roughly as in the last years. Compared to last year, the support for outputting proofs was extended, and should now cover all relevant theory modules for CASC (but not yet all proof strategies).

7.16 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 [49]. 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

Prover9 has available many strategies; the following statements apply to CASC.

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.

For CASC Prover9 uses KBO to order terms for demodulation and for the inference rules, with a simple rule for determining symbol precedence.

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 [48]. (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 [74] is used for forward and backward nonunit subsumption. Prover9 is available from

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

Expected Competition Performance

Prover9 is the CASC fixed point, against which progress can be judged. Each year it is expected do worse than the previous year, relative to the other systems.

7.17 Satallax 3.2

Michael Färber
Universität Innsbruck, Austria

Architecture

Satallax 3.2 [14] 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 [22] is responsible for much of the proof search. The theoretical basis of search is a complete ground tableau calculus for higher-order logic [16] with a choice operator [3]. Problems are given in the THF format.

Proof search: 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 [15]. These formulae and propositional clauses correspond to instances of the tableau rules. Satallax uses the SAT solver MiniSat to test the current set of propositional clauses for unsatisfiability. If the clauses are unsatisfiable, then the original branch is unsatisfiable. Optionally, Satallax generates first-order formulae in addition to the propositional clauses. If this option is used, then Satallax periodically calls the first-order theorem prover E [76] to test for first-order unsatisfiability. If the set of first-order formulae is unsatisfiable, then the original branch is unsatisfiable. Upon request, Satallax attempts to reconstruct a proof which can be output in the TSTP format. The proof reconstruction has been significantly changed since Satallax 3.0 in order to make proof reconstruction more efficient and thus less likely to fail within the time constraints.

Strategies

There are about 150 flags that control the order in which formulae and instantiation terms are considered and propositional clauses are generated. Other flags activate some optional extensions to the basic proof procedure (such as whether or not to call the theorem prover E). A collection of flag settings is called a mode. Approximately 500 modes have been defined and tested so far. 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. Before deciding on the schedule to use, Satallax parses the problem and determines if it is big enough that a SInE-based premise selection algorithm [30] should be used. If SInE is not activated, then Satallax 3.2 uses a strategy schedule consisting of 37 modes. Each mode is tried for time limits ranging from less than a second to just over 1 minute. If SInE is activated, then Satallax is run with a SInE-specific schedule consisting of 20 possible SInE parameter values selecting different premises and some corresponding modes and time limits.

Implementation

Satallax is implemented in OCaml, making use of the external tools MiniSat (via a foreign function interface) and E. Satallax is available at:

<http://satallaxprover.com>

Expected Competition Performance

Satallax 3.2 is the CASC-26 THF division winner.

7.18 Satallax 3.3

Michael Färber
Universität Innsbruck, Austria

Architecture

Satallax 3.3 [14] 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 [22] is responsible for much of the proof search. The theoretical basis of search is a complete ground tableau calculus for higher-order logic [16] with a choice operator [3]. Problems are given in the THF format.

Proof search: 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 [15]. These formulae and propositional clauses correspond to instances of the tableau rules. Satallax uses the SAT solver MiniSat to test the current set of propositional clauses for unsatisfiability. If the clauses are unsatisfiable, then the original branch is unsatisfiable. Optionally, Satallax generates first-order formulae in addition to the propositional clauses. If this option is used, then Satallax periodically calls the first-order theorem prover E [76] to test for first-order unsatisfiability. If the set of first-order formulae is unsatisfiable, then the original branch is unsatisfiable. Upon request, Satallax attempts to reconstruct a proof which can be output in the TSTP format.

Strategies

There are about 150 flags that control the order in which formulae and instantiation terms are considered and propositional clauses are generated. Other flags activate some optional extensions to the basic proof procedure (such as whether or not to call the theorem prover E). A collection of flag settings is called a mode. Approximately 500 modes have been defined and tested so far. 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. Before deciding on the schedule to use, Satallax parses the problem and determines if it is big enough that a SInE-based premise selection algorithm [30] should be used. If SInE is not activated, then Satallax 3.3 uses a strategy schedule consisting of 48 modes (16 of which make use of Mizar style soft types). Each mode is tried for time limits ranging from less than a second to just over 1 minute. If SInE is activated, then Satallax is run with a SInE-specific schedule consisting of 20 possible SInE parameter values selecting different premises and some corresponding modes and time limits.

Implementation

Satallax is implemented in OCaml, making use of the external tools MiniSat (via a foreign function interface) and E. Satallax is available at:

<http://satallaxprover.com>

Expected Competition Performance

Satallax 3.3 adds support for Mizar style soft types. As this technique has only minor impact on proof search performance, we expect Satallax 3.3 to perform about as well as Satallax 3.2.

7.19 Twee 2.2

Nick Smallbone
Chalmers University of Technology, Sweden

Architecture

Twee 2.2 is an equational prover based on unifying completion [2]. It features ground joinability testing [47] and a connectedness test [1], which together eliminate many redundant inferences in the presence of unorientable equations.

Twee’s implementation of ground joinability testing performs case splits on the order of variables, in the style of [47], and discharges individual cases by rewriting modulo a variable ordering. It is able to pick only useful case splits and to case split on a subset of the variables, which makes it efficient enough to be switched on unconditionally.

Horn clauses are encoded as equations as described in [18]. The CASC version of Twee “handles” non-Horn clauses by discarding them.

Strategies

Twee’s strategy is simple and it does not tune its heuristics or strategy based on the input problem. The term ordering is always KBO; functions are ordered by number of occurrences (commonly-occurring symbols are smaller) and always have weight 1.

The main loop is a DISCOUNT loop. The active set contains rewrite rules and unorientable equations, which are used for rewriting, and the passive set contains unprocessed critical pairs. Twee often interreduces the active set, and occasionally simplifies the passive set with respect to the active set. Each critical pair is scored using a weighted sum of the weight of both of its terms. Terms are treated as DAGs when computing weights, i.e., duplicate subterms are only counted once per term. The weights of critical pairs that correspond to Horn clauses are adjusted by the heuristic described in [18], section 5.

Implementation

Twee is written in Haskell. Terms are represented as array-based flatterms for efficient unification and matching. Rewriting uses an imperfect discrimination tree.

The passive set is represented as a heap. It achieves high space efficiency (12 bytes per critical pair) by storing the parent rule numbers and overlap position instead of the full critical pair and by grouping all critical pairs of each rule into one heap entry.

Twee uses an LCF-style kernel: all rules in the active set come with a certified proof object which traces back to the input axioms. When a conjecture is proved, the proof object is transformed into a human-readable proof. Proof construction does not harm efficiency because the proof kernel is invoked only when a new rule is accepted. In particular, reasoning about

the passive set does not invoke the kernel. The translation from Horn clauses to equations is not yet certified.

Twee can be downloaded from:

<http://nick8325.github.io/twee>

Expected Competition Performance

Twee is an equational prover dressed up as a Horn clause prover, so its performance on general first-order problems will be poor. It will do well on Horn problems with a heavy equational component. With a bit of luck, it may solve a few hard problems.

7.20 Vampire 4.0

Giles Reger

University of Manchester, United Kingdom

Architecture

Vampire 4.0 is an automatic theorem prover for first-order logic. Vampire implements the calculi of ordered binary resolution and superposition for handling equality. It also implements the Inst-gen calculus and a MACE-style finite model builder. Splitting in resolution-based proof search is controlled by the AVATAR architecture, which uses a SAT solver to make splitting decisions. Both resolution and instantiation based proof search make use of global subsumption.

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. Internally, Vampire works only with clausal normal form. Problems in the full first-order logic syntax are clausified during preprocessing. Vampire implements many useful preprocessing transformations including the Sine axiom selection algorithm.

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

Strategies

Vampire 4.0 provides a very large number of options for strategy selection. The most important ones are:

- Choices of saturation algorithm:
 - Limited Resource Strategy
 - DISCOUNT loop
 - Otter loop
 - Instantiation using the Inst-Gen calculus
 - MACE-style finite model building with sort inference
- Splitting via AVATAR

- 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.
- Ground equational reasoning via congruence closure.
- Evaluation of interpreted functions.
- Extensionality resolution with detection of extensionality axioms

Implementation

Vampire 4.0 is implemented in C++.

Expected Competition Performance

Vampire 4.0 is the CASC-26 LTB division winner.

7.21 Vampire 4.1

Giles Regeer
University of Manchester, United Kingdom

Architecture

Vampire [43] 4.1 is an automatic theorem prover for first-order logic. Vampire implements the calculi of ordered binary resolution and superposition for handling equality. It also implements the Inst-gen calculus [40] and a MACE-style finite model builder [62]. Splitting in resolution-based proof search is controlled by the AVATAR architecture [131] which uses a SAT or SMT solver to make splitting decisions. Both resolution and instantiation based proof search make use of global subsumption [40].

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. Internally, Vampire works only with clausal normal form. Problems in the full first-order logic syntax are clausified during preprocessing. Vampire implements many useful preprocessing transformations including the SinE axiom selection algorithm.

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

Strategies

Vampire 4.1 provides a very large number of options for strategy selection. The most important ones are:

- Choices of saturation algorithm:
 - Limited Resource Strategy [69].
 - DISCOUNT loop
 - Otter loop
 - Instantiation using the Inst-Gen calculus
 - MACE-style finite model building with sort inference
- Splitting via AVATAR
- 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.
- Ground equational reasoning via congruence closure.
- Addition of theory axioms and evaluation of interpreted functions.
- Use of Z3 [19] with AVATAR to restrict search to ground-theory-consistent splitting branches.
- Extensionality resolution [26] with detection of extensionality axioms.

Implementation

Vampire 4.1 is implemented in C++.

Expected Competition Performance

Vampire 4.1 is the CASC-26 TFA and FNT division winner.

7.22 Vampire 4.2

Giles Reger
University of Manchester, United Kingdom

Architecture

Vampire [43] 4.2 is an automatic theorem prover for first-order logic. Vampire implements the calculi of ordered binary resolution and superposition for handling equality. It also implements the Inst-gen calculus and a MACE-style finite model builder [62]. Splitting in resolution-based proof search is controlled by the AVATAR architecture which uses a SAT or SMT solver to make splitting decisions [131, 61]. Both resolution and instantiation based proof search make use of global subsumption.

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. Internally, Vampire works only with clausal normal form. Problems in the full first-order logic syntax are clausified during preprocessing. Vampire implements many useful preprocessing transformations including the SinE axiom selection algorithm. When a theorem is proved, the system produces a verifiable proof, which validates both the clausification phase and the refutation of the CNF.

Strategies

Vampire 4.2 provides a very large number of options for strategy selection. The most important ones are:

- Choices of saturation algorithm:
 - Limited Resource Strategy [69]
 - DISCOUNT loop
 - Otter loop
 - Instantiation using the Inst-Gen calculus
 - MACE-style finite model building with sort inference
- Splitting via AVATAR [131]
- A variety of optional simplifications.
- Parameterized reduction orderings.
- A number of built-in literal selection functions and different modes of comparing literals [29].
- Age-weight ratio that specifies how strongly lighter clauses are preferred for inference selection.
- Set-of-support strategy.
- Ground equational reasoning via congruence closure.

- Addition of theory axioms and evaluation of interpreted functions.
- Use of Z3 with AVATAR to restrict search to ground-theory-consistent splitting branches [61].
- Specialised theory instantiation and unification
- Extensionality resolution with detection of extensionality axioms

Implementation

Vampire 4.2 is implemented in C++. It makes use of minisat and z3.

Expected Competition Performance

Vampire 4.2 is the CASC-26 FOF division winner.

7.23 Vampire 4.3

Giles Reger

University of Manchester, United Kingdom

This description is very similar to that of Vampire 4.2. The main difference is the use of theory instantiation and unification with abstraction [63] for theory reasoning (this was experimental in 4.2). The set-of-support strategy for theory reasoning has also been extended. Little has changed in other areas of Vampire. As always there have been some small improvements to heuristics, data structures and schedules but nothing fundamentally new.

Architecture

Vampire [43] 4.3 is an automatic theorem prover for first-order logic. Vampire implements the calculi of ordered binary resolution and superposition for handling equality. It also implements the Inst-gen calculus and a MACE-style finite model builder [62]. Splitting in resolution-based proof search is controlled by the AVATAR architecture which uses a SAT or SMT solver to make splitting decisions [131, 61]. Both resolution and instantiation based proof search make use of global subsumption.

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. Internally, Vampire works only with clausal normal form. Problems in the full first-order logic syntax are clausified during preprocessing. Vampire implements many useful preprocessing transformations including the SinE axiom selection algorithm. When a theorem is proved, the system produces a verifiable proof, which validates both the clausification phase and the refutation of the CNF.

Strategies

Vampire 4.3 provides a very large number of options for strategy selection. The most important ones are:

- Choices of saturation algorithm:
 - Limited Resource Strategy [69]
 - DISCOUNT loop
 - Otter loop
 - Instantiation using the Inst-Gen calculus
 - MACE-style finite model building with sort inference
- Splitting via AVATAR [131]
- A variety of optional simplifications.
- Parameterized reduction orderings.
- A number of built-in literal selection functions and different modes of comparing literals [29].
- Age-weight ratio that specifies how strongly lighter clauses are preferred for inference selection.
- Set-of-support strategy.
- Ground equational reasoning via congruence closure.
- Addition of theory axioms and evaluation of interpreted functions.
- Use of Z3 with AVATAR to restrict search to ground-theory-consistent splitting branches [61].
- Specialised theory instantiation and unification [63]
- Extensionality resolution with detection of extensionality axioms

Implementation

Vampire 4.3 is implemented in C++. It makes use of minisat and z3.

Expected Competition Performance

Vampire 4.3 should be better than 4.2 at theory reasoning and as good as 4.2 at everything else.

8 Conclusion

The 9th IJCAR ATP System Competition was the twenty-third 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 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] L. Bachmair and N. Dershowitz. Critical Pair Criteria for Completion. *Journal of Symbolic Computation*, 6(1):1–18, 1988.
- [2] L. Bachmair, N. Dershowitz, and D.A. Plaisted. Completion Without Failure. In H. Ait-Kaci and M. Nivat, editors, *Resolution of Equations in Algebraic Structures*, pages 1–30. Academic Press, 1989.
- [3] J. Backes and C.E. Brown. Analytic Tableaux for Higher-Order Logic with Choice. *Journal of Automated Reasoning*, 47(4):451–479, 2011.
- [4] H. Barbosa, P. Fontaine, and A. Reynolds. Congruence Closure with Free Variables. In A. Legay and T. Margaria, editors, *Proceedings of the 23rd International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, number 10205 in Lecture Notes in Computer Science, pages 2134–230. Springer-Verlag, 2017.
- [5] C. Barrett, C. Conway, M. Deters, L. Hadarean, D. Jovanovic, T. King, A. Reynolds, and C. Tinelli. CVC4. In G. Gopalakrishnan and S. Qadeer, editors, *Proceedings of the 23rd International Conference on Computer Aided Verification*, number 6806 in Lecture Notes in Computer Science, pages 171–177. Springer-Verlag, 2011.
- [6] 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.
- [7] 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.
- [8] 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 G. Dowek, 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.
- [9] C. Benzmüller, F. Rabe, and G. Sutcliffe. THF0 - The Core TPTP Language for Classical Higher-Order Logic. In P. Baumgartner, A. Armando, and G. Dowek, 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.
- [10] C. Benzmüller, V. Sorge, M. Jamnik, and M. Kerber. Combined Reasoning by Automated Cooperation. *Journal of Applied Logic*, 6(3):318–342, 2008.
- [11] W. Bibel. *Automated Theorem Proving*. Vieweg and Sohn, 1987.
- [12] A. Biere. PicoSAT Essentials. *Journal on Satisfiability, Boolean Modeling and Computation*, 4:75–97, 2008.
- [13] J. Blanchette and A. Paskevich. TFF1: The TPTP Typed First-order Form with Rank-1 Polymorphism. In M.P. Bonacina, editor, *Proceedings of the 24th International Conference on Auto-*

- ated Deduction*, number 7898 in Lecture Notes in Artificial Intelligence, pages 414–420. Springer-Verlag, 2013.
- [14] 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*, number 7364 in Lecture Notes in Artificial Intelligence, pages 111–117, 2012.
 - [15] C.E. Brown. Reducing Higher-Order Theorem Proving to a Sequence of SAT Problems. *Journal of Automated Reasoning*, 51(1):57–77, 2013.
 - [16] 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.
 - [17] K. Claessen, A. Lilliestrom, and N. Smallbone. Sort It Out with Monotonicity - Translating between Many-Sorted and Unsorted First-Order Logic. In N. Bjørner 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.
 - [18] K. Claessen and N. Smallbone. Efficient Encodings of First-Order Horn Formulas in Equational Logic. In D. Galmiche, S. Schulz, and R. Sebastiani, editors, *Proceedings of the 9th International Joint Conference on Automated Reasoning*, page To appear, 2018.
 - [19] L. de Moura and N. Bjørner. 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.
 - [20] H. de Nivelle. Classical Logic with Partial Functions. *Journal of Automated Reasoning*, 47(4):399–425, 2011.
 - [21] H. de Nivelle. Theorem Proving for Classical Logic with Partial Functions by Reduction to Kleene Logic. *Journal of Logic and Computation*, 27(2):509–548, 2017.
 - [22] 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.
 - [23] 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.
 - [24] 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.
 - [25] 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.
 - [26] A. Gupta, L. Kovacs, B. Kragl, and A. Voronkov. Extensional Crisis and Proving Identity. In F. Cassez and J-F. Franck, editors, *Proceedings of the 12th International Symposium on Automated Technology for Verification and Analysis*, number 8837 in Lecture Notes in Computer Science, pages 185–200, 2014.
 - [27] J. Hernandez and K. Korovin. Towards an Abstraction-Refinement Framework for Reasoning with Large Theories. In T. Eiter, D. Sands, S. Schulz, J. Urban, G. Sutcliffe, and A. Voronkov, editors, *Proceedings of the IWIL Workshop and LPAR Short Presentations*, number 1 in Kalpa Publications in Computing, 2017.
 - [28] K. Hoder, Z. Khasidashvili, K. Korovin, and A. Voronkov. Preprocessing Techniques for First-Order Clausification. In G. Cabodi and S. Singh, editors, *Proceedings of the Formal Methods in Computer-Aided Design 2012*, pages 44–51. IEEE Press, 2012.

- [29] K. Hoder, G. Regeer, M. Suda, and A. Voronkov. Selecting the Selection. In N. Olivetti and A. Tiwari, editors, *Proceedings of the 8th International Joint Conference on Automated Reasoning*, number 9706 in Lecture Notes in Artificial Intelligence, pages 313–329, 2016.
- [30] 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.
- [31] F. Hutter, H. Hoos, K. Leyton-Brown, and T. Stützle. ParamILS: An Automatic Algorithm Configuration Framework. *Journal of Artificial Intelligence Research*, 36:267–306, 2009.
- [32] J. Jakubuv, M. Suda, and J. Urban. Automated Invention of Strategies and Term Orderings for Vampire. In C. Benzmüller, C. Lisetti, and M. Theobald, editors, *Proceedings of the 3rd Global Conference on Artificial Intelligence*, number 50 in EPiC Series in Computing, pages 121–133. EasyChair Publications, 2017.
- [33] J. Jakubuv and J. Urban. BliStrTune: Hierarchical Invention of Theorem Proving Strategies. In Y. Bertot and V. Vafeiadis, editors, *Proceedings of Certified Programs and Proofs 2017*, pages 43–52. ACM, 2017.
- [34] J. Jakubuv and J. Urban. Hierarchical Invention of Theorem Proving Strategies. *AI Communications*, 31(3):237–250, 2018.
- [35] C. Kaliszyk, G. Sutcliffe, and F. Rabe. TH1: The TPTP Typed Higher-Order Form with Rank-1 Polymorphism. In P. Fontaine, S. Schulz, and J. Urban, editors, *Proceedings of the 5th Workshop on the Practical Aspects of Automated Reasoning*, number 1635 in CEUR Workshop Proceedings, pages 41–55, 2016.
- [36] C. Kaliszyk and J. Urban. Stronger Automation for Flyspeck by Feature Weighting and Strategy Evolution. In *Proceedings of the 3rd International Workshop on Proof Exchange for Theorem Proving*, page To appear. EasyChair Proceedings in Computing, 2013.
- [37] C. Kaliszyk, J. Urban, and J. Vyskocil. Machine Learner for Automated Reasoning 0.4 and 0.5. In S. Schulz, L. de Moura, and B. Konev, editors, *Proceedings of the 4th Workshop on Practical Aspects of Automated Reasoning*, number 31 in EPiC Series in Computing, pages 60–66, 2015.
- [38] Z. Khasidashvili and K. Korovin. Predicate Elimination for Preprocessing in First-order Theorem Proving. In N. Creignou and D. Le Berre, editors, *Proceedings of the 19th International Conference on Theory and Applications of Satisfiability Testing*, number 9710 in Lecture Notes in Computer Science, pages 361–372. Springer-Verlag, 2016.
- [39] K. Korovin. iProver - An Instantiation-Based Theorem Prover for First-order Logic (System Description). In P. Baumgartner, A. Armando, and G. Dowek, editors, *Proceedings of the 4th International Joint Conference on Automated Reasoning*, number 5195 in Lecture Notes in Artificial Intelligence, pages 292–298, 2008.
- [40] K. Korovin. Inst-Gen - A Modular Approach to Instantiation-Based Automated Reasoning. In A. Voronkov and C. Weidenbach, editors, *Programming Logics, Essays in Memory of Harald Ganzinger*, number 7797 in Lecture Notes in Computer Science, pages 239–270. Springer-Verlag, 2013.
- [41] K. Korovin. Non-cyclic Sorts for First-order Satisfiability. In P. Fontaine, C. Ringeissen, and R. Schmidt, editors, *Proceedings of the International Symposium on Frontiers of Combining Systems*, number 8152 in Lecture Notes in Computer Science, pages 214–228, 2013.
- [42] K. Korovin and C. Stickse. A Note on Model Representation and Proof Extraction in the First-order Instantiation-based Calculus Inst-Gen. In R. Schmidt and F. Papacchini, editors, *Proceedings of the 19th Automated Reasoning Workshop*, pages 11–12, 2012.
- [43] L. Kovacs and A. Voronkov. First-Order Theorem Proving and Vampire. In N. Sharygina and H. Veith, editors, *Proceedings of the 25th International Conference on Computer Aided Verification*, number 8044 in Lecture Notes in Artificial Intelligence, pages 1–35. Springer-Verlag, 2013.

- [44] 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.
- [45] B. Loechner. Things to Know When Implementing KBO. *Journal of Automated Reasoning*, 36(4):289–310, 2006.
- [46] B. Loechner. Things to Know When Implementing LBO. *Journal of Artificial Intelligence Tools*, 15(1):53–80, 2006.
- [47] U. Martin and T. Nipkow. Ordered Rewriting and Confluence. In M.E. Stickel, editor, *Proceedings of the 10th International Conference on Automated Deduction*, number 449 in Lecture Notes in Artificial Intelligence, pages 366–380. Springer-Verlag, 1990.
- [48] W.W. McCune. Solution of the Robbins Problem. *Journal of Automated Reasoning*, 19(3):263–276, 1997.
- [49] W.W. McCune. Otter 3.3 Reference Manual. Technical Report ANL/MS-C-263, Argonne National Laboratory, Argonne, USA, 2003.
- [50] 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.
- [51] 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 G. Dowek, editors, *Proceedings of the 4th International Joint Conference on Automated Reasoning*, number 5195 in Lecture Notes in Artificial Intelligence, pages 283–291, 2008.
- [52] J. Otten. Restricting Backtracking in Connection Calculi. *AI Communications*, 23(2-3):159–182, 2010.
- [53] J. Otten. A Non-clausal Connection Calculus. In K. Brunnler and G. Metcalfe, editors, *Proceedings of the 20th International Conference on Automated Reasoning with Analytic Tableaux and Related Methods*, number 6793 in Lecture Notes in Artificial Intelligence, pages 226–241. Springer-Verlag, 2011.
- [54] J. Otten. Implementing Connection Calculi for First-order Modal Logics. In K. Korovin and S. Schulz, editors, *Proceedings of the 9th International Workshop on the Implementation of Logics*, number 22 in EPiC Series in Computing, pages 18–32, 2012.
- [55] J. Otten. MleanCoP: A Connection Prover for First-Order Modal Logic. In S. Demri, D. Kapur, and C. Weidenbach, editors, *Proceedings of the 7th International Joint Conference on Automated Reasoning*, number 8562 in Lecture Notes in Artificial Intelligence, pages 269–276, 2014.
- [56] J. Otten. nanoCoP: A Non-clausal Connection Prover. In S. Demri, D. Kapur, and C. Weidenbach, editors, *Proceedings of the 7th International Joint Conference on Automated Reasoning*, number 8562 in Lecture Notes in Artificial Intelligence, pages 302–312, 2016.
- [57] J. Otten. nanoCoP: Natural Non-clausal Theorem Proving. In C. Sierra, editor, *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 4924–4928. Morgan-Kaufmann, 2017.
- [58] J. Otten. Non-clausal Connection Calculi for Non-classical Logics. In C. Nalon and R. Schmidt, editors, *Proceedings of the 26th International Conference on Automated Reasoning with Analytic Tableaux and Related Methods*, number 10501 in Lecture Notes in Artificial Intelligence, pages 209–227. Springer-Verlag, 2017.
- [59] J. Otten. Proof Search Optimizations for Non-clausal Connection Calculi. In B. Konev, P. Rümmer, and J. Urban, editors, *Proceedings of the 6th Workshop on Practical Aspects of Automated Reasoning*, 2018.
- [60] J. Otten and W. Bibel. leanCoP: Lean Connection-Based Theorem Proving. *Journal of Symbolic Computation*, 36(1-2):139–161, 2003.
- [61] G. Reger, N. Bjørner, M. Suda, and A. Voronkov. AVATAR Modulo Theories. In C. Benz Müller,

- G. Sutcliffe, and R. Rojas, editors, *Proceedings of the 2nd Global Conference on Artificial Intelligence*, number 41 in EPiC Series in Computing, pages 39–52. EasyChair Publications, 2016.
- [62] G. Reger, M. Suda, and A. Voronkov. Finding Finite Models in Multi-Sorted First Order Logic. In N. Creignou and D. Le Berre, editors, *Proceedings of the 19th International Conference on Theory and Applications of Satisfiability Testing*, number 9710 in Lecture Notes in Computer Science, pages 323–341. Springer-Verlag, 2016.
- [63] G. Reger, M. Suda, and A. Voronkov. Unification with Abstraction and Theory Instantiation in Saturation-Based Reasoning. In D. Beyer and M. Huisman, editors, *Proceedings of the 24th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, number 10805 in Lecture Notes in Computer Science, pages 3–22. Springer-Verlag, 2018.
- [64] A. Reynolds, M. Deters, V. Kuncak, C. Barrett, and C. Tinelli. Counterexample Guided Quantifier Instantiation for Synthesis in CVC4. In D. Kroening and C. Pasareanu, editors, *Proceedings of the 27th International Conference on Computer Aided Verification*, number 9207 in Lecture Notes in Computer Science, pages 198–216. Springer-Verlag, 2015.
- [65] A. Reynolds, C. Tinelli, and L. de Moura. Finding Conflicting Instances of Quantified Formulas in SMT. In K. Claessen and V. Kuncak, editors, *Proceedings of the 14th Conference on Formal Methods in Computer-Aided Design*, pages 195–202, 2014.
- [66] A. Reynolds, C. Tinelli, A. Goel, and S. Krstic. Finite Model Finding in SMT. In N. Sharygina and H. Veith, editors, *Proceedings of the 25th International Conference on Computer Aided Verification*, number 8044 in Lecture Notes in Computer Science, pages 640–655. Springer-Verlag, 2013.
- [67] A. Reynolds, C. Tinelli, A. Goel, S. Krstic, M. Deters, and C. Barrett. Quantifier Instantiation Techniques for Finite Model Finding in SMT. In M.P. Bonacina, editor, *Proceedings of the 24th International Conference on Automated Deduction*, number 7898 in Lecture Notes in Artificial Intelligence, pages 377–391. Springer-Verlag, 2013.
- [68] A. Reynolds, C. Tinelli, D. Jovanovic, and C. Barrett. Designing Theory Solvers with Extensions. In C. Dixon and M. Finger, editors, *Proceedings of the 11th International Symposium on Frontiers of Combining Systems*, number 10483 in Lecture Notes in Computer Science, pages 22–40. Springer-Verlag, 2017.
- [69] A. Riazanov and A. Voronkov. Limited Resource Strategy in Resolution Theorem Proving. *Journal of Symbolic Computation*, 36(1-2):101–115, 2003.
- [70] 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.
- [71] P. Rümmer. E-Matching with Free Variables. In N. Bjørner 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 359–374. Springer-Verlag, 2012.
- [72] 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.
- [73] S. Schulz. E: A Brainiac Theorem Prover. *AI Communications*, 15(2-3):111–126, 2002.
- [74] 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. Springer-Verlag, 2004.
- [75] 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*, number 7364 in Lecture Notes in Artificial Intelligence, pages 477–483. Springer-Verlag, 2012.

- [76] S. Schulz. Simple and Efficient Clause Subsumption with Feature Vector Indexing. In M.P. Bonacina and M. Stickel, editors, *Automated Reasoning and Mathematics: Essays in Memory of William W. McCune*, number 7788 in Lecture Notes in Artificial Intelligence, pages 45–67. Springer-Verlag, 2013.
- [77] S. Schulz. System Description: E 1.8. In K. McMillan, A. Middeldorp, and A. Voronkov, editors, *Proceedings of the 19th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning*, number 8312 in Lecture Notes in Computer Science, pages 477–483. Springer-Verlag, 2013.
- [78] A. Steen and C. Benzmüller. The Higher-Order Prover Leo-III. In D. Galmiche, S. Schulz, and R. Sebastiani, editors, *Proceedings of the 8th International Joint Conference on Automated Reasoning*, Lecture Notes in Artificial Intelligence, page To appear, 2018.
- [79] G. Sutcliffe. Proceedings of the CADE-16 ATP System Competition. Trento, Italy, 1999.
- [80] G. Sutcliffe. Proceedings of the CADE-17 ATP System Competition. Pittsburgh, USA, 2000.
- [81] G. Sutcliffe. The CADE-16 ATP System Competition. *Journal of Automated Reasoning*, 24(3):371–396, 2000.
- [82] G. Sutcliffe. Proceedings of the IJCAR ATP System Competition. Siena, Italy, 2001.
- [83] G. Sutcliffe. The CADE-17 ATP System Competition. *Journal of Automated Reasoning*, 27(3):227–250, 2001.
- [84] G. Sutcliffe. Proceedings of the CADE-18 ATP System Competition. Copenhagen, Denmark, 2002.
- [85] G. Sutcliffe. Proceedings of the CADE-19 ATP System Competition. Miami, USA, 2003.
- [86] G. Sutcliffe. Proceedings of the 2nd IJCAR ATP System Competition. Cork, Ireland, 2004.
- [87] G. Sutcliffe. Proceedings of the CADE-20 ATP System Competition. Tallinn, Estonia, 2005.
- [88] G. Sutcliffe. The IJCAR-2004 Automated Theorem Proving Competition. *AI Communications*, 18(1):33–40, 2005.
- [89] G. Sutcliffe. Proceedings of the 3rd IJCAR ATP System Competition. Seattle, USA, 2006.
- [90] G. Sutcliffe. The CADE-20 Automated Theorem Proving Competition. *AI Communications*, 19(2):173–181, 2006.
- [91] G. Sutcliffe. Proceedings of the CADE-21 ATP System Competition. Bremen, Germany, 2007.
- [92] G. Sutcliffe. The 3rd IJCAR Automated Theorem Proving Competition. *AI Communications*, 20(2):117–126, 2007.
- [93] G. Sutcliffe. Proceedings of the 4th IJCAR ATP System Competition. Sydney, Australia, 2008.
- [94] G. Sutcliffe. The CADE-21 Automated Theorem Proving System Competition. *AI Communications*, 21(1):71–82, 2008.
- [95] 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.
- [96] G. Sutcliffe. Proceedings of the CADE-22 ATP System Competition. Montreal, Canada, 2009.
- [97] G. Sutcliffe. The 4th IJCAR Automated Theorem Proving System Competition - CASC-J4. *AI Communications*, 22(1):59–72, 2009.
- [98] G. Sutcliffe. Proceedings of the 5th IJCAR ATP System Competition. Edinburgh, United Kingdom, 2010.
- [99] G. Sutcliffe. The CADE-22 Automated Theorem Proving System Competition - CASC-22. *AI Communications*, 23(1):47–60, 2010.
- [100] G. Sutcliffe. Proceedings of the CADE-23 ATP System Competition. Wroclaw, Poland, 2011.
- [101] G. Sutcliffe. The 5th IJCAR Automated Theorem Proving System Competition - CASC-J5. *AI Communications*, 24(1):75–89, 2011.

- [102] G. Sutcliffe. Proceedings of the 6th IJCAR ATP System Competition. Manchester, England, 2012.
- [103] G. Sutcliffe. The CADE-23 Automated Theorem Proving System Competition - CASC-23. *AI Communications*, 25(1):49–63, 2012.
- [104] G. Sutcliffe. Proceedings of the 24th CADE ATP System Competition. Lake Placid, USA, 2013.
- [105] G. Sutcliffe. The 6th IJCAR Automated Theorem Proving System Competition - CASC-J6. *AI Communications*, 26(2):211–223, 2013.
- [106] G. Sutcliffe. Proceedings of the 7th IJCAR ATP System Competition. Vienna, Austria, 2014.
- [107] G. Sutcliffe. The CADE-24 Automated Theorem Proving System Competition - CASC-24. *AI Communications*, 27(4):405–416, 2014.
- [108] G. Sutcliffe. Proceedings of the CADE-25 ATP System Competition. Berlin, Germany, 2015. <http://www.tptp.org/CASC/25/Proceedings.pdf>.
- [109] G. Sutcliffe. The 7th IJCAR Automated Theorem Proving System Competition - CASC-J7. *AI Communications*, 28(4):683–692, 2015.
- [110] G. Sutcliffe. Proceedings of the 8th IJCAR ATP System Competition. Coimbra, Portugal, 2016. <http://www.tptp.org/CASC/J8/Proceedings.pdf>.
- [111] G. Sutcliffe. The 8th IJCAR Automated Theorem Proving System Competition - CASC-J8. *AI Communications*, 29(5):607–619, 2016.
- [112] G. Sutcliffe. Proceedings of the 26th CADE ATP System Competition. Gothenburg, Sweden, 2017. <http://www.tptp.org/CASC/26/Proceedings.pdf>.
- [113] G. Sutcliffe. The CADE-26 Automated Theorem Proving System Competition - CASC-26. *AI Communications*, 30(6):419–432, 2017.
- [114] G. Sutcliffe. The TPTP Problem Library and Associated Infrastructure. From CNF to TH0, TPTP v6.4.0. *Journal of Automated Reasoning*, 59(4):483–502, 2017.
- [115] G. Sutcliffe and C. Benz Müller. Automated Reasoning in Higher-Order Logic using the TPTP THF Infrastructure. *Journal of Formalized Reasoning*, 3(1):1–27, 2010.
- [116] 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.
- [117] 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.
- [118] G. Sutcliffe and C. Suttner. The CADE-18 ATP System Competition. *Journal of Automated Reasoning*, 31(1):23–32, 2003.
- [119] G. Sutcliffe and C. Suttner. The CADE-19 ATP System Competition. *AI Communications*, 17(3):103–182, 2004.
- [120] G. Sutcliffe, C. Suttner, and F.J. Pelletier. The IJCAR ATP System Competition. *Journal of Automated Reasoning*, 28(3):307–320, 2002.
- [121] G. Sutcliffe and C.B. Suttner. Special Issue: The CADE-13 ATP System Competition. *Journal of Automated Reasoning*, 18(2), 1997.
- [122] G. Sutcliffe and C.B. Suttner. The Procedures of the CADE-13 ATP System Competition. *Journal of Automated Reasoning*, 18(2):163–169, 1997.
- [123] G. Sutcliffe and C.B. Suttner. Proceedings of the CADE-15 ATP System Competition. Lindau, Germany, 1998.
- [124] G. Sutcliffe and C.B. Suttner. The CADE-15 ATP System Competition. *Journal of Automated Reasoning*, 23(1):1–23, 1999.
- [125] G. Sutcliffe and C.B. Suttner. Evaluating General Purpose Automated Theorem Proving Systems. *Artificial Intelligence*, 131(1-2):39–54, 2001.

- [126] G. Sutcliffe and J. Urban. The CADE-25 Automated Theorem Proving System Competition - CASC-25. *AI Communications*, 29(3):423–433, 2016.
- [127] C.B. Suttner and G. Sutcliffe. The CADE-14 ATP System Competition. *Journal of Automated Reasoning*, 21(1):99–134, 1998.
- [128] 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.
- [129] J. Urban. BliStr: The Blind Strategymaker. arXiv:1301.2683, 2013.
- [130] 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 G. Dowek, 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.
- [131] A. Voronkov. AVATAR: The New Architecture for First-Order Theorem Provers. In A. Biere and R. Bloem, editors, *Proceedings of the 26th International Conference on Computer Aided Verification*, number 8559 in Lecture Notes in Computer Science, pages 696–710, 2014.
- [132] M. Wisniewski, A. Steen, and C. Benzmüller. LeoPARD - A Generic Platform for the Implementation of Higher-Order Reasoners. In M. Kerber, J. Carette, C. Kaliszyk, F. Rabe, and V. Sorge, editors, *Proceedings of the International Conference on Intelligent Computer Mathematics*, number 9150 in Lecture Notes in Computer Science, pages 325–330. Springer-Verlag, 2015.
- [133] Y. Xu, S. Liu, J. Chen, X. Zhong, and X. He. Contradiction Separation Based Dynamic Multi-clause Synergized Automated Deduction. *Information Sciences*, 462:93–113, 2018.