

# Hints for AVATAR (and some more)

Martin Suda

Czech Technical University in Prague, Czech Republic

PIWo 2019, Prague, October 2019

**Some people actually use ATPs to do math!**

## Some people actually use ATPs to do math!

- e.g., Bob Veroff and Michael Kinyon
- using Otter, Prover9, Mace4
- questions from algebra: axioms bases for boolean algebras, ortho-lattices, loop theory
- targeting open problems (e.g. the AIM conjecture)

## Some people actually use ATPs to do math!

- e.g., Bob Veroff and Michael Kinyon
- using Otter, Prover9, Mace4
- questions from algebra: axioms bases for boolean algebras, ortho-lattices, loop theory
- targeting open problems (e.g. the AIM conjecture)

## In what sense interactive?

- a single proof attempt (ATP call) usually does not solve it
- trying different formulations / axiomatizations
- trying various additional assumptions and learning from them

## Some people actually use ATPs to do math!

- e.g., Bob Veroff and Michael Kinyon
- using Otter, Prover9, Mace4
- questions from algebra: axioms bases for boolean algebras, ortho-lattices, loop theory
- targeting open problems (e.g. the AIM conjecture)

## In what sense interactive?

- a single proof attempt (ATP call) usually does not solve it
- trying different formulations / axiomatizations
- trying various additional assumptions and learning from them

➡ By the way, these attempts may run for weeks!

## What is a hint?

- a clause supplied by the user as part of the input
- whenever a newly derived clause  $C$  subsumes a hint clause, this  $C$  is prioritized for selection

## What is a hint?

- a clause supplied by the user as part of the input
  - whenever a newly derived clause  $C$  subsumes a hint clause, this  $C$  is prioritized for selection
- ➔ Hints are a means for steering the proof search!

## What is a hint?

- a clause supplied by the user as part of the input
- whenever a newly derived clause  $C$  subsumes a hint clause, this  $C$  is prioritized for selection

➔ Hints are a means for steering the proof search!

## Where do hints come from?

- the (expert) user just thinks of some



## What is a hint?

- a clause supplied by the user as part of the input
- whenever a newly derived clause  $C$  subsumes a hint clause, this  $C$  is prioritized for selection

➔ Hints are a means for steering the proof search!

## Where do hints come from?

- the (expert) user just thinks of some
- more realistically: clauses from proofs of similar theorems or of the same theorem but under different assumptions

## What is a hint?

- a clause supplied by the user as part of the input
- whenever a newly derived clause  $C$  subsumes a hint clause, this  $C$  is prioritized for selection

➔ Hints are a means for steering the proof search!

## Where do hints come from?

- the (expert) user just thinks of some
- more realistically: clauses from proofs of similar theorems or of the same theorem but under different assumptions

➔ Hope that similar theorems can be proved using similar intermediate steps.

## What is a hint?

- a clause supplied by the user as part of the input
- whenever a newly derived clause  $C$  subsumes a hint clause, this  $C$  is prioritized for selection

➡ Hints are a means for steering the proof search!

## Where do hints come from?

- the (expert) user just thinks of some
- more realistically: clauses from proofs of similar theorems or of the same theorem but under different assumptions

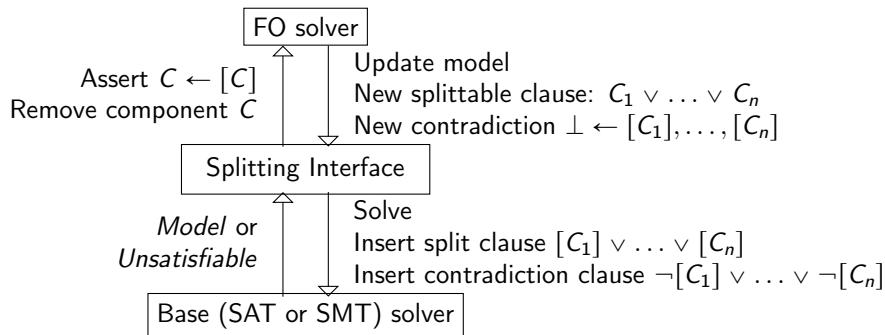
➡ Hope that similar theorems can be proved using similar intermediate steps.

## How to come up with hints automatically?

## AVATAR [Voronkov'14]

- modern architecture of first order theorem provers
- integrates saturation with a SAT solver (or an SMT solver)
- efficient realization of the *clause splitting rule*
- instead of one monolithic proof search  
a sequence of proof searches on (much) smaller sub-problems
  
- implemented in theorem prover Vampire
- shown highly successful in practice

# AVATAR architecture overview



Instead of waiting for the user to supply hints for problem  $P$  ...

... attempt  $P$  using AVATAR and collect as hints the first-order parts of the clauses appearing in the sub-proofs of the so far derived contradiction clauses

Instead of waiting for the user to supply hints for problem  $P$  ...

... attempt  $P$  using AVATAR and collect as hints the first-order parts of the clauses appearing in the sub-proofs of the so far derived contradiction clauses

**DEMO!**

- 1 Hints for AVATAR
- 2 An Experiment
- 3 What is a Significant Improvement?



- 1 Hints for AVATAR
- 2 An Experiment
- 3 What is a Significant Improvement?

## Vampire setup:

- `--saturation_algorithm discount` (for stability)
- `--age_weight_ratio 1:10` (works well with discount)
- `--time_limit 10` (reasonable time to finish)

## Vampire setup:

- `--saturation_algorithm discount` (for stability)
- `--age_weight_ratio 1:10` (works well with discount)
- `--time_limit 10` (reasonable time to finish)

## Computers:

- either Starexec
- or CTU's (slurm) cluster

## Vampire setup:

- `--saturation_algorithm discount` (for stability)
- `--age_weight_ratio 1:10` (works well with discount)
- `--time_limit 10` (reasonable time to finish)

## Computers:

- either Starexec
- or CTU's (slurm) cluster

## The benchmark:

- TPTP v 7.2.0
- 17573 eligible first-order problems

(on Starexec)

configuration	solved	uniques	additional
base	7914	0	7914
base+hints	7882	2	62
sac	8100	13	299
sac+hints	8106	13	23

base = -sa discount -awr 10 -t 10

sac = --split\_at\_activation on

(on Starexec)

configuration	solved	uniques	additional
base	7914	0	7914
base+hints	7882	2	62
sac	8100	13	299
sac+hints	8106	13	23

base = -sa discount -awr 10 -t 10

sac = --split\_at\_activation on

Experimented with AVATAR flushing; also not very interesting

## MIZAR bushy “small”

- 57 880 problems translated from the MIZAR library

## MIZAR bushy “small”

- 57 880 problems translated from the MIZAR library

(base: -sa discount -awr 10 -t 10 -sac on)



## MIZAR bushy “small”

- 57 880 problems translated from the MIZAR library

(base: `-sa discount -awr 10 -t 10 -sac on`)

## Results

configuration	solved	uniques
base	14843	184
base+hints	14873	214

## MIZAR bushy “small”

- 57 880 problems translated from the MIZAR library

(base: `-sa discount -awr 10 -t 10 -sac on`)

## Results

configuration	solved	uniques
base	14843	184
base+hints	14873	214

(30 problems is approx. 0.5‰ of the benchmark size)

So, should we be sad and abandon the idea?

So, should we be sad and abandon the idea?

**Maybe, but ...**

**Maybe, but ...**

- maybe it only gets interesting with really hard problems!

# So, should we be sad and abandon the idea?

## Maybe, but . . .

- maybe it only gets interesting with really hard problems!
- maybe we should have a smarter notion of similarity!
  - demodulate hints?

## Maybe, but . . .

- maybe it only gets interesting with really hard problems!
- maybe we should have a smarter notion of similarity!
  - demodulate hints?
- maybe we need restarts to prevent the prover from choking

## Maybe, but . . .

- maybe it only gets interesting with really hard problems!
- maybe we should have a smarter notion of similarity!
  - demodulate hints?
- maybe we need restarts to prevent the prover from choking
  
- we should also try strengthening the theory with reasonable additional assumptions, as routinely done by Veroff et al.



## Maybe, but . . .

- maybe it only gets interesting with really hard problems!
- maybe we should have a smarter notion of similarity!
  - demodulate hints?
- maybe we need restarts to prevent the prover from choking
  
- we should also try strengthening the theory with reasonable additional assumptions, as routinely done by Veroff et al.

➡ Ongoing and future work!

- 1 Hints for AVATAR
- 2 An Experiment
- 3 What is a Significant Improvement?

**When should we get excited about a new technique?**

**When should we get excited about a new technique?**

- 1 The idea looks clever and sophisticated

## When should we get excited about a new technique?

- 1 The idea looks clever and sophisticated
  - ➡ Could aim for a pure theory paper at CADE!

## When should we get excited about a new technique?

- 1 The idea looks clever and sophisticated
  - ➡ Could aim for a pure theory paper at CADE!
- 2 Solves more problems than baseline

## When should we get excited about a new technique?

- 1 The idea looks clever and sophisticated
  - ➔ Could aim for a pure theory paper at CADE!
- 2 Solves more problems than baseline
  - ➔ Obviously, this gives us more power!

## When should we get excited about a new technique?

- 1 The idea looks clever and sophisticated
  - ➡ Could aim for a pure theory paper at CADE!
- 2 Solves more problems than baseline
  - ➡ Obviously, this gives us more power!
- 3 The solution set differs enough from baseline



## When should we get excited about a new technique?

- 1 The idea looks clever and sophisticated
  - ➔ Could aim for a pure theory paper at CADE!
- 2 Solves more problems than baseline
  - ➔ Obviously, this gives us more power!
- 3 The solution set differs enough from baseline
  - ➔ To have a chance to improve strategy schedule ...

Proof search is known to be very fragile. Even a small change will “stir” it and create a different solution set.

Proof search is known to be very fragile. Even a small change will “stir” it and create a different solution set.

**What does it mean to differ enough from baseline?**

Proof search is known to be very fragile. Even a small change will “stir” it and create a different solution set.

## What does it mean to differ enough from baseline?

- Keep a database of problems known to be solvable by some strategy and compare against that.

Proof search is known to be very fragile. Even a small change will “stir” it and create a different solution set.

## What does it mean to differ enough from baseline?

- Keep a database of problems known to be solvable by some strategy and compare against that.
  - ➔ Computationally expensive, open ended, but YES!

Proof search is known to be very fragile. Even a small change will “stir” it and create a different solution set.

## What does it mean to differ enough from baseline?

- Keep a database of problems known to be solvable by some strategy and compare against that.
  - ➔ Computationally expensive, open ended, but YES!

While the database is being built . . .

Proof search is known to be very fragile. Even a small change will “stir” it and create a different solution set.

## What does it mean to differ enough from baseline?

- Keep a database of problems known to be solvable by some strategy and compare against that.
  - ➡ Computationally expensive, open ended, but YES!

While the database is being built . . .

Let's have some random fun!

# Randomly permuting the input problem

Use `tptp4X -trandomize` from the TPTP toolset to:

- randomize the order of commutative logical operations
- randomize the order of formulas



# Randomly permuting the input problem

Use `tptp4X -trandomize` from the TPTP toolset to:

- randomize the order of commutative logical operations
- randomize the order of formulas

Can we solve more problems?

Use `tptp4X -trandomize` from the **TPTP** toolset to:

- randomize the order of commutative logical operations
- randomize the order of formulas

**Can we solve more problems?**

configuration	solved	uniques	additional
straight	8612	53	8612
shuffled1	8773	60	345
shuffled2	8788	85	128
shuffled3	8775	48	48

(now on the CTU cluster)

Use `tptp4X -trandomize` from the **TPTP** toolset to:

- randomize the order of commutative logical operations
- randomize the order of formulas

Can we solve more problems?

configuration	solved	uniques	additional
straight	8612	53	8612
shuffled1	8773	60	345
shuffled2	8788	85	128
shuffled3	8775	48	48

(now on the CTU cluster)

➡ recalling **randoCoP** (Raths, Otten; 2008)

## Clause Selection and Age-weight Ratio

Vampire alternates between selecting the next given clause by age (old first) and by weight (light first) under a given ratio.

## Clause Selection and Age-weight Ratio

Vampire alternates between selecting the next given clause by age (old first) and by weight (light first) under a given ratio.

Normally, this alternation is regular. What if we change it to probabilistic?

## Clause Selection and Age-weight Ratio

Vampire alternates between selecting the next given clause by age (old first) and by weight (light first) under a given ratio.

Normally, this alternation is regular. What if we change it to probabilistic?

configuration	solved	uniques	additional
base	8725	12	8725
rnd1	8747	8	91
rnd2	8744	16	37
rnd3	8768	23	37
rnd4	8735	14	21
rnd5	8741	16	16

```
base = -sa discount -awr 1:1 -t 10
```

## Empirical research with an ATP:

- 1 have a new idea
- 2 implement (and debug)
- 3 conduct experiments

## Empirical research with an ATP:

- 1 have a new idea
- 2 implement (and debug)
- 3 conduct experiments

## When are the results significant?

- improving overall performance (high total solved)
- solving hard problems (“the uniques”)



## Empirical research with an ATP:

- 1 have a new idea
- 2 implement (and debug)
- 3 conduct experiments

## When are the results significant?

- improving overall performance (high total solved)
- solving hard problems (“the uniques”)

Why don't we use (carefully seeded) randomness to prove more theorems (without much actual extra thinking)?