

# Two Learning Operators for Clause Selection Guidance: An Experimental Evaluation

Martin Suda\*

Czech Technical University in Prague, Czech Republic

AITP, Aussios, September 2024

# Machine Learning Boosted Automated Theorem Proving

**ATP technology:**

**ATP technology:** saturation-based

**ATP technology:** saturation-based

- state of the art (cf. CASC)
- E, iProver, Vampire, . . .

**ATP technology:** saturation-based

- state of the art (cf. CASC)
- E, iProver, Vampire, . . .

**Heuristic to boost:**

**ATP technology:** saturation-based

- state of the art (cf. CASC)
- E, iProver, Vampire, . . .

**Heuristic to boost:** clause selection

**ATP technology:** saturation-based

- state of the art (cf. CASC)
- E, iProver, Vampire, . . .

**Heuristic to boost:** clause selection

- the most important choice point
- “selecting the proof clauses” intuition



**ATP technology:** saturation-based

- state of the art (cf. CASC)
- E, iProver, Vampire, . . .

**Heuristic to boost:** clause selection

- the most important choice point
- “selecting the proof clauses” intuition

**Two main approaches to date:**

## **ATP technology:** saturation-based

- state of the art (cf. CASC)
- E, iProver, Vampire, . . .

## **Heuristic to boost:** clause selection

- the most important choice point
- “selecting the proof clauses” intuition

## **Two main approaches to date:**

- ENIGMA-style
- RL-inspired

**ATP technology:** saturation-based

- state of the art (cf. CASC)
- E, iProver, Vampire, . . .

**Heuristic to boost:** clause selection

- the most important choice point
- “selecting the proof clauses” intuition

**Two main approaches to date:**

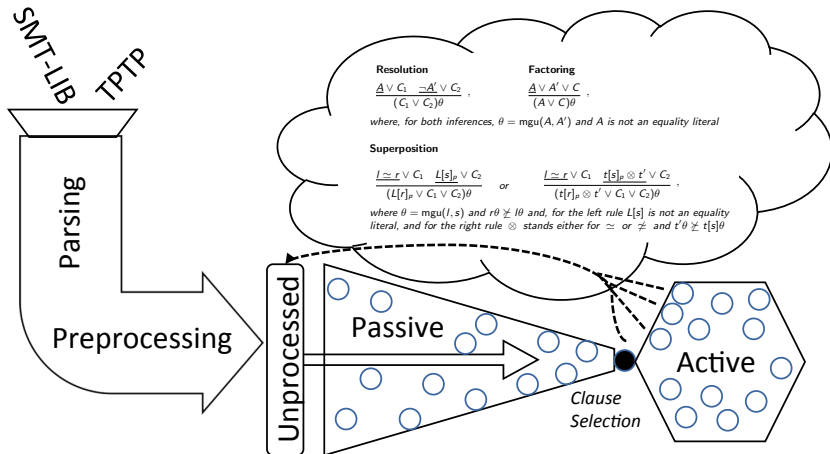
- ENIGMA-style
- RL-inspired

What are the differences? What is the same? Which one is better?

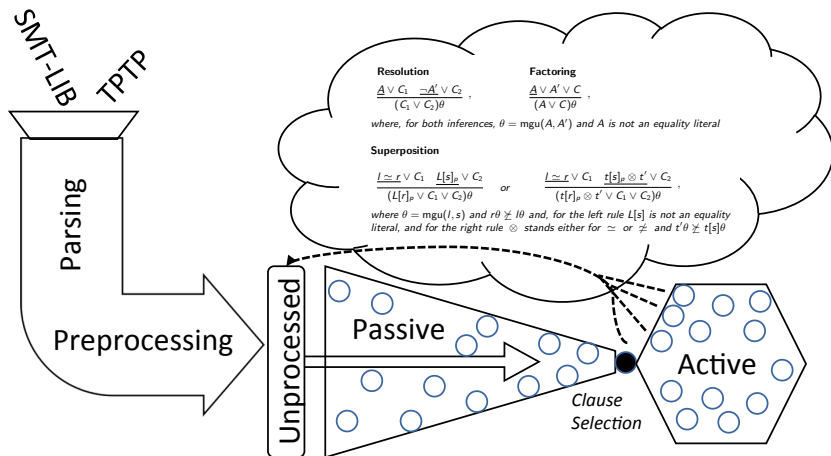
- 1 Saturation and Clause Selection
- 2 ENIGMA-style Guidance
- 3 RL-Inspired Guidance
- 4 Compare, Contrast, Evaluate

- 1 Saturation and Clause Selection
- 2 ENIGMA-style Guidance
- 3 RL-Inspired Guidance
- 4 Compare, Contrast, Evaluate

# Saturation-based Theorem Proving



# Saturation-based Theorem Proving



At a typical successful end:  $|Passive| \gg |Active| \gg |Proof|$

# How is clause selection traditionally done?

## Take simple clause evaluation criteria:

- age: prefer clauses that were generated long time ago
- weight: prefer clauses with fewer symbols



# How is clause selection traditionally done?

## Take simple clause evaluation criteria:

- age: prefer clauses that were generated long time ago
- weight: prefer clauses with fewer symbols

## Combine them into a single scheme:

- have a priority queue ordering *Passive* for each criterion
- alternate between selecting from the queues using a fixed ratio

# How is clause selection traditionally done?

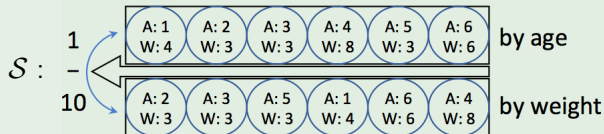
## Take simple clause evaluation criteria:

- age: prefer clauses that were generated long time ago
- weight: prefer clauses with fewer symbols

## Combine them into a single scheme:

- have a priority queue ordering *Passive* for each criterion
- alternate between selecting from the queues using a fixed ratio

## Example (Organizing *Passive* via two priority queues)



- 1 Saturation and Clause Selection
- 2 ENIGMA-style Guidance**
- 3 RL-Inspired Guidance
- 4 Compare, Contrast, Evaluate

## The core idea

*Learn to recognize and prefer for selection clauses that look like those that contributed to a proof in past successful runs.*

➔ [Schulz00], ENIGMA [Jakubův&Urban17], [Loos et al.'17], ...

## The core idea

*Learn to recognize and prefer for selection clauses that look like those that contributed to a proof in past successful runs.*

➔ [Schulz00], ENIGMA [Jakubův&Urban17], [Loos et al.'17], ...

## The “pos/neg”s of E:

E prover can be asked to output, for every clause selected in a run, whether it ended up in the final proof (**pos**) or not (**neg**)

## The core idea

*Learn to recognize and prefer for selection clauses that look like those that contributed to a proof in past successful runs.*

➔ [Schulz00], ENIGMA [Jakubův&Urban17], [Loos et al.'17], ...

## The “pos/neg”s of E:

E prover can be asked to output, for every clause selected in a run, whether it ended up in the final proof (**pos**) or not (**neg**)

## Next comes the ML:

- represent those clauses somehow (features, NNs, ...)
- train a binary classifier on the task
- integrate back with the prover:

## The core idea

*Learn to recognize and prefer for selection clauses that look like those that contributed to a proof in past successful runs.*

➔ [Schulz00], ENIGMA [Jakubův&Urban17], [Loos et al.'17], ...

## The “pos/neg”s of E:

E prover can be asked to output, for every clause selected in a run, whether it ended up in the final proof (**pos**) or not (**neg**)

## Next comes the ML:

- represent those clauses somehow (features, NNs, ...)
- train a binary classifier on the task
- integrate back with the prover: “try to do more of the **pos**”



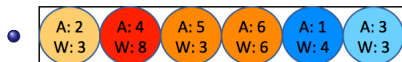




# Possible Ways of Integrating the Learnt Advice

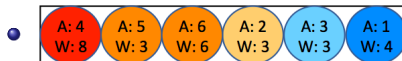
## Priority:

- sort by model's Y/N and tiebreak by age

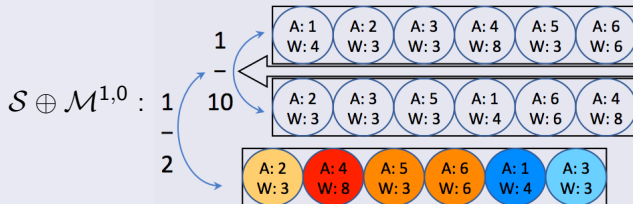


## Logits:

- even a binary classifier internally uses a real value



## Combine with the original strategy



- 1 Saturation and Clause Selection
- 2 ENIGMA-style Guidance
- 3 RL-Inspired Guidance**
- 4 Compare, Contrast, Evaluate

# Saturation as an Reinforcement-Learning Environment

What ATP heuristics would the aliens come up with?

What ATP heuristics would the aliens come up with?

## Agent

- the clause selection heuristic

## Action

- the next clause to select from the current passive set

What ATP heuristics would the aliens come up with?

## Agent

- the clause selection heuristic

## Action

- the next clause to select from the current passive set

## State

- static - the conjecture we are trying to prove
- evolving - the internal state of the prover at particular moment

What ATP heuristics would the aliens come up with?

## Agent

- the clause selection heuristic

## Action

- the next clause to select from the current passive set

## State

- static - the conjecture we are trying to prove
- evolving - the internal state of the prover at particular moment

## Reward

- Score 1 point for solving a problem (within the time limit)

What ATP heuristics would the aliens come up with?

## Agent

- the clause selection heuristic

## Action

- the next clause to select from the current passive set

## State

- static - the conjecture we are trying to prove
- evolving - the internal state of the prover at particular moment

## Reward

- Score 1 point for solving a problem (within the time limit) ???



What ATP heuristics would the aliens come up with?

## Agent

- the clause selection heuristic

## Action

- the next clause to select from the current passive set

## State

- static - the conjecture we are trying to prove
- evolving - the internal state of the prover at particular moment

## Reward

- Score 1 point for solving a problem (within the time limit) ???

➡ TRAIL [Crouse et al.'21], [McKeown'23], [Shminke'23], ...

# Policy Gradient and REINFORCE [Williams'92]

The (evolving) state  $s$  of an ATP is a large amorphous blob:

- value-based methods (Q-learning, DQN, ...) seem hopeless

# Policy Gradient and REINFORCE [Williams'92]

The (evolving) state  $s$  of an ATP is a large amorphous blob:

- value-based methods (Q-learning, DQN, ...) seem hopeless

Instead, with policy gradient methods, we train a network to directly predict the policy  $\pi(a|s; \theta)$

# Policy Gradient and REINFORCE [Williams'92]

The (evolving) state  $s$  of an ATP is a large amorphous blob:

- value-based methods (Q-learning, DQN, ...) seem hopeless

Instead, with policy gradient methods, we train a network to directly predict the policy  $\pi(a|s; \theta)$

- can sample actions according to the distribution  $\pi$

# Policy Gradient and REINFORCE [Williams'92]

The (evolving) state  $s$  of an ATP is a large amorphous blob:

- value-based methods (Q-learning, DQN, ...) seem hopeless

Instead, with policy gradient methods, we train a network to directly predict the policy  $\pi(a|s; \theta)$

- can sample actions according to the distribution  $\pi$
- imperfect information  $\Rightarrow$  the optimal policy may be stochastic!

# Policy Gradient and REINFORCE [Williams'92]

The (evolving) state  $s$  of an ATP is a large amorphous blob:

- value-based methods (Q-learning, DQN, ...) seem hopeless

Instead, with policy gradient methods, we train a network to directly predict the policy  $\pi(a|s; \theta)$

- can sample actions according to the distribution  $\pi$
- imperfect information  $\Rightarrow$  the optimal policy may be stochastic!

## Policy Gradient Theorem

$$\nabla_{\theta} v_{\pi}(s_{initial}) \propto \mathbb{E}_{s \sim \mu} \mathbb{E}_{a \sim \pi} q_{\pi}(s, a) \nabla_{\theta} \ln \pi(a|s; \theta)$$

The (evolving) state  $s$  of an ATP is a large amorphous blob:

- value-based methods (Q-learning, DQN, ...) seem hopeless

Instead, with policy gradient methods, we train a network to directly predict the policy  $\pi(a|s; \theta)$

- can sample actions according to the distribution  $\pi$
- imperfect information  $\Rightarrow$  the optimal policy may be stochastic!

## Policy Gradient Theorem

$$\nabla_{\theta} v_{\pi}(s_{initial}) \propto \mathbb{E}_{s \sim \mu} \mathbb{E}_{a \sim \pi} q_{\pi}(s, a) \nabla_{\theta} \ln \pi(a|s; \theta)$$

**The devil in the details:**

- with  $\pi(C|s_i; \theta) = \text{softmax}([\text{NN}_{\theta}(\text{features}_C)]_{C \in \text{Passive}_i})$ , the “ $\nabla_{\theta} \ln \pi$ ”-bit boils down to the usual NLL loss

# Policy Gradient and REINFORCE [Williams'92]

The (evolving) state  $s$  of an ATP is a large amorphous blob:

- value-based methods (Q-learning, DQN, ...) seem hopeless

Instead, with policy gradient methods, we train a network to directly predict the policy  $\pi(a|s; \theta)$

- can sample actions according to the distribution  $\pi$
- imperfect information  $\Rightarrow$  the optimal policy may be stochastic!

## Policy Gradient Theorem

$$\nabla_{\theta} v_{\pi}(s_{initial}) \propto \mathbb{E}_{s \sim \mu} \mathbb{E}_{a \sim \pi} q_{\pi}(s, a) \nabla_{\theta} \ln \pi(a|s; \theta)$$

**The devil in the details:**

- with  $\pi(C|s_i; \theta) = \text{softmax}([\text{NN}_{\theta}(\text{features}_C)]_{C \in \text{Passive}_i})$ , the “ $\nabla_{\theta} \ln \pi$ ”-bit boils down to the usual NLL loss
- for  $q_{\pi}(s, C)$  we simply pick **Did C show up in the found proof?**



- 1 Saturation and Clause Selection
- 2 ENIGMA-style Guidance
- 3 RL-Inspired Guidance
- 4 Compare, Contrast, Evaluate**

## Starts with:

- ENIGMA-style: a working clause selection heuristic
- RL-inspired: “tabula rasa”

## Starts with:

- ENIGMA-style: a working clause selection heuristic
- RL-inspired: “tabula rasa”

## Training data:

- ENIGMA-style: pos/neg; over selected only (static)
- RL-inspired: traces; over all the generated (changes in time)

## Starts with:

- ENIGMA-style: a working clause selection heuristic
- RL-inspired: “tabula rasa”

## Training data:

- ENIGMA-style: pos/neg; over selected only (static)
- RL-inspired: traces; over all the generated (changes in time)

## Attractor:

- Both: clauses from found proofs

## Starts with:

- ENIGMA-style: a working clause selection heuristic
- RL-inspired: “tabula rasa”

## Training data:

- ENIGMA-style: pos/neg; over selected only (static)
- RL-inspired: traces; over all the generated (changes in time)

## Attractor:

- Both: clauses from found proofs

## Integrating the learned advice:

- ENIGMA-style: combine with your original heuristic
- RL-inspired: one queue sorted by the predicted scores

## Model:

- ENIGMA-style: a binary classifier
- RL-inspired: regression (logits)  $\Rightarrow$  action probabilities

## Model:

- ENIGMA-style: a binary classifier
- RL-inspired: regression (logits)  $\Rightarrow$  action probabilities

## Loss function (for the neural incarnations):

- ENIGMA-style: binary cross entropy (NLL)
- RL-inspired: weighted NLL (weights  $\sim$  returns)

## Model:

- ENIGMA-style: a binary classifier
- RL-inspired: regression (logits)  $\Rightarrow$  action probabilities

## Loss function (for the neural incarnations):

- ENIGMA-style: binary cross entropy (NLL)
- RL-inspired: weighted NLL (weights  $\sim$  returns)

## Iterative improvement:

- Both: yes (ENIGMA calls it “looping”)



## Architecture

- simple clause features: age, weight, pos/neg-length, justEq/justNeq, varOcc, goalDist, numSplits
- a neural part:  $\text{MLP}(\text{features}_C) \rightarrow \text{logit}$

## Architecture

- simple clause features: age, weight, pos/neg-length, justEq/justNeq, varOcc, goalDist, numSplits
- a neural part:  $\text{MLP}(\text{features}_C) \rightarrow \text{logit}$

## Experimental setup

- extend Vampire theorem prover
- 3000 randomly select TPTP problems (FOF/CNF)
- time limit:  $\sim 10$ s per problem

## Architecture

- simple clause features: age, weight, pos/neg-length, justEq/justNeq, varOcc, goalDist, numSplits
- a neural part:  $\text{MLP}(\text{features}_C) \rightarrow \text{logit}$

## Experimental setup

- extend Vampire theorem prover
- 3000 randomly select TPTP problems (FOF/CNF)
- time limit:  $\sim 10\text{s}$  per problem

## ENIGMA-style

- unfortunately, did not manage to beat the baseline

## Architecture

- simple clause features: age, weight, pos/neg-length, justEq/justNeq, varOcc, goalDist, numSplits
- a neural part:  $\text{MLP}(\text{features}_C) \rightarrow \text{logit}$

## Experimental setup

- extend Vampire theorem prover
- 3000 randomly select TPTP problems (FOF/CNF)
- time limit:  $\sim 10$ s per problem

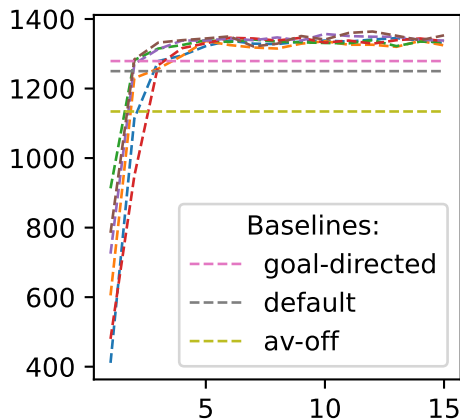
## ENIGMA-style

- unfortunately, did not manage to beat the baseline

## RL-inspired

- can beat the default strategy by 6%,  
(a good goal-directed strategy by 3.5%) on the test set

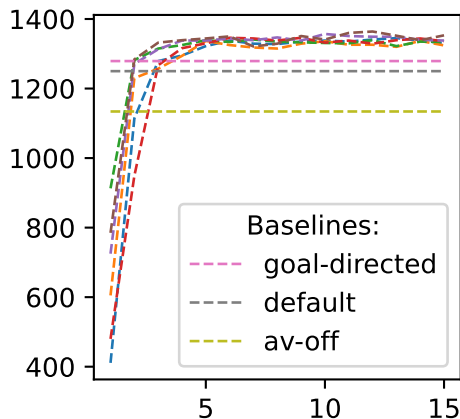
# Training with the RL-inspired Operator



## Observations:

- random initialization  
⇒ different start performance

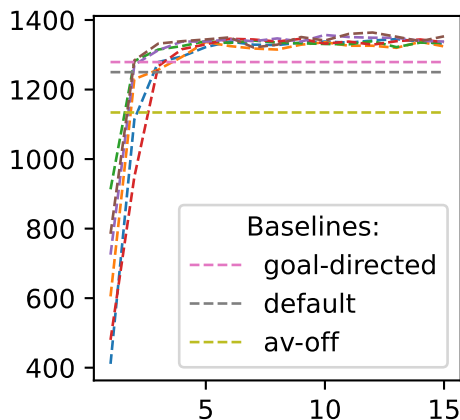
# Training with the RL-inspired Operator



## Observations:

- random initialization  
⇒ different start performance
- hidden layer size 16, 32, 64, 128, 256 all similar performance

# Training with the RL-inspired Operator

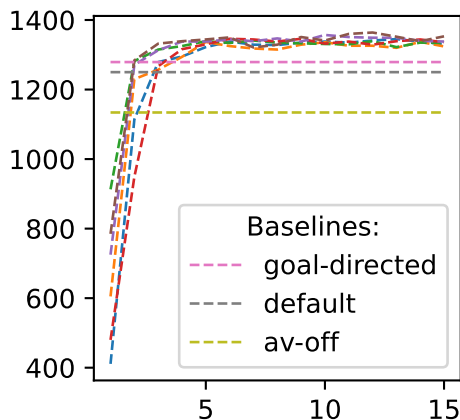


## Observations:

- random initialization  
⇒ different start performance
- hidden layer size 16, 32, 64, 128, 256 all similar performance

A first (modest) ML-based improvement of a SoTA ATP on TPTP!

# Training with the RL-inspired Operator



## Observations:

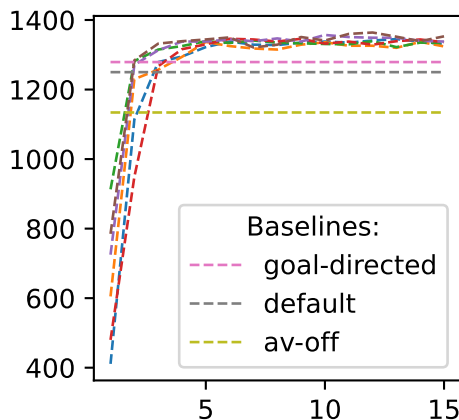
- random initialization  
⇒ different start performance
- hidden layer size 16, 32, 64, 128, 256 all similar performance

A first (modest) ML-based improvement of a SoTA ATP on TPTP!

Future work: ready for experiments with richer feature sets!



# Training with the RL-inspired Operator



## Observations:

- random initialization  
⇒ different start performance
- hidden layer size 16, 32, 64, 128, 256 all similar performance

A first (modest) ML-based improvement of a SoTA ATP on TPTP!

Future work: ready for experiments with richer feature sets!

Thank you!