

COMBINING LEARNING AND REASONING OVER LARGE FORMAL MATH CORPORA

Josef Urban

Czech Technical University in Prague

Dagstuhl, September 3rd, 2019



European Research Council
Established by the European Commission

Outline

Motivation, Learning vs. Reasoning

Computer Understandable (Formal) Math

Learning of Theorem Proving - Overview

Demos

High-level Reasoning Guidance: Premise Selection

Low Level Guidance of Theorem Provers

Mid-level Reasoning Guidance

Autoformalization

Motivation: Learning vs. Reasoning

“C’est par la logique qu’on démontre, c’est par l’intuition qu’on invente.”

(It is by logic that we prove, but by intuition that we discover.)

Henri Poincaré, Mathematical Definitions and Education.

“Hypothesen sind Netze; nur der fängt, wer auswirft.”

(Hypotheses are nets: only he who casts will catch.)

Novalis, quoted by Popper – The Logic of Scientific Discovery

How Do We Automate Math and Science?

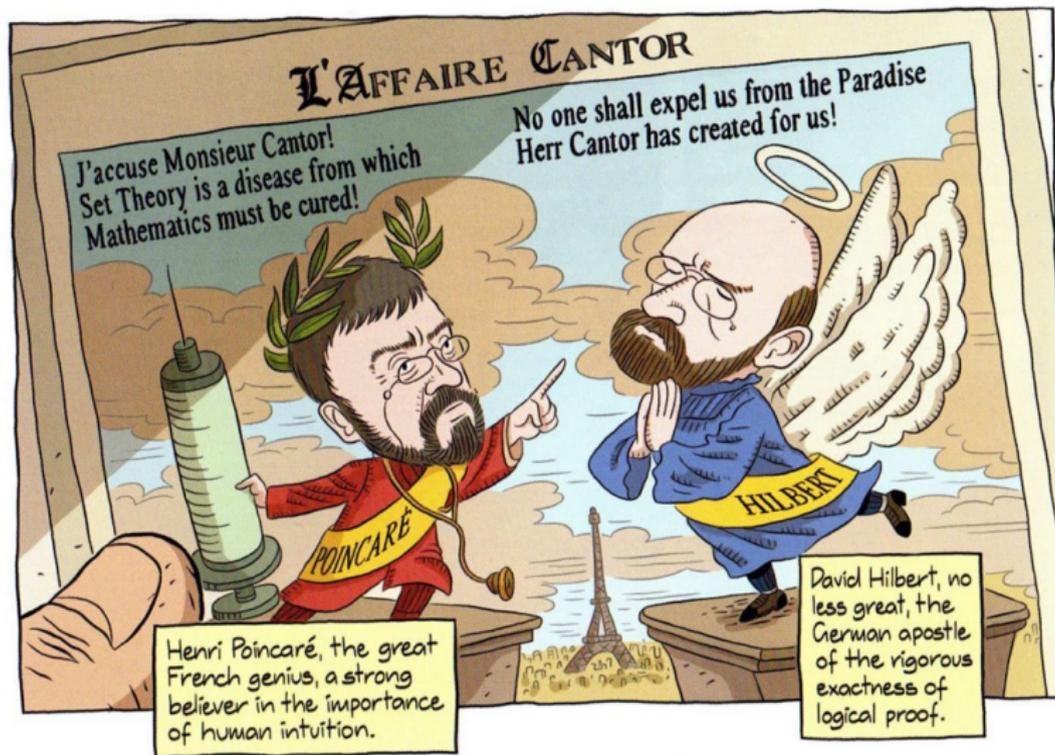
- What is mathematical and scientific thinking?
- Pattern-matching, analogy, induction from examples
- Deductive reasoning
- Complicated feedback loops between induction and deduction
- Using a lot of previous knowledge - both for induction and deduction

- We need to develop such methods on computers
- Are there any large corpora suitable for nontrivial deduction?
- Yes! Large libraries of formal proofs and theories
- So let's develop strong AI on them!

History, Motivation, AI/TP/ML

- Intuition vs Formal Reasoning – Poincaré vs Hilbert, Science & Method
- Turing's 1950 paper: **Learning Machines**, learn Chess?, undecidability??
- 50s-60s: Beginnings of ATP and ITP – Davis, Simon, Robinson, de Bruijn
- Lenat, Langley: **AM**, manually-written heuristics, learn Kepler laws,...
- Denzinger, Schulz, Goller, Fuchs – late 90's, ATP-focused:
Learning from Previous Proof Experience
- My MSc (1998): Try ILP to learn rules and heuristics from IMPS/Mizar
- Since: Use large formal math (Big Proof) corpora: Mizar, Isabelle, HOL ... to combine/develop symbolic/statistical deductive/inductive ML/TP/AI ... hammer-style methods, internal guidance, **feedback loops**, ...
- AI vs ML vs DL?: Ben Goertzel's 2018 Prague talk:
<https://youtu.be/Zt2HSTuGBn8>

Intuition vs Formal Reasoning – Poincaré vs Hilbert



[Adapted from: *Logicomix: An Epic Search for Truth* by A. Doxiadis]

Induction/Learning vs Reasoning – Henri Poincaré



- Science and Method: Ideas about the interplay between correct deduction and induction/intuition
- *“And in demonstration itself logic is not all. The true **mathematical reasoning is a real induction** [...]”*
- I believe he was right: strong general reasoning engines have to **combine deduction and induction** (learning patterns from data, making conjectures, etc.)

Learning vs Reasoning – Alan Turing 1950 – AI



- 1950: *Computing machinery and intelligence* – AI, Turing test
- “We may hope that machines will eventually compete with men in *all purely intellectual fields*.” (regardless of his 1936 undecidability result!)
- last section on **Learning Machines**:
- “But which are the best ones [fields] to start [learning on] with?”
- “... Even this is a difficult decision. Many people think that a very abstract activity, like the *playing of chess*, would be best.”
- Why not try with **math**? It is much more (universally?) expressive ...

Why Combine Learning and Reasoning Today?

1 It practically helps!

- Automated theorem proving for large formal verification is **useful**:
 - Formal Proof of the Kepler Conjecture (2014 – Hales – 20k lemmas)
 - Formal Proof of the Feit-Thompson Theorem (2012 – Gonthier)
 - Verification of compilers (CompCert) and microkernels (seL4)
 - ...
- **But** good learning/AI methods needed to cope with large theories!

2 Blue Sky AI Visions:

- Get **strong AI** by learning/reasoning over large KBs of **human thought**?
- Big formal theories: good **semantic** approximation of such thinking KBs?
- Deep non-contradictory semantics – better than scanning books?
- Gradually try **learning math/science**:
 - What are the components (inductive/deductive thinking)?
 - How to combine them together?

The AITP Plan for World Domination

- 1 Make **large formal thought** (Mizar/MML, Isabelle/HOL/AFP, HOL/Flyspeck ...) accessible to strong reasoning and learning AI tools – **DONE** (or well under way)
- 2 Test/Use/Evolve existing AI and ATP tools on such large corpora
- 3 Build custom/combined inductive/deductive tools/metasystems
- 4 Continuously test performance, define harder AI tasks as the performance grows

Hilbert's update for 21st century (AGI'18):

NO ONE SHALL DRIVE US FROM THE SEMANTIC AI PARADISE
OF COMPUTER-UNDERSTANDABLE MATH AND SCIENCE!

aitp-conference.org - since 2016

Outline

Motivation, Learning vs. Reasoning

Computer Understandable (Formal) Math

Learning of Theorem Proving - Overview

Demos

High-level Reasoning Guidance: Premise Selection

Low Level Guidance of Theorem Provers

Mid-level Reasoning Guidance

Autoformalization

What is Formal Mathematics?

- Developed thanks to the Leibniz/Russell/Frege/Hilbert/... program
- Mathematics put on **formal logic foundations** (*symbolic computation*)
- ... which btw. led also to the rise of computers (Turing/Church, 1930s)
- Formal math (1950/60s): **combine formal foundations and computers**
- For AGI/Singularity people: Formal proof is the *Secure Hardware Environment* from Vinge's Rainbows End
- **Conceptually very simple:**
- Write all your axioms and theorems so that computer understands them
- Write all your inference rules so that computer understands them
- Use the computer to check that your proofs follow the rules
- **But in practice, it turns out not to be so simple**
- Many approaches, still not mainstream, but big breakthroughs recently

The QED Manifesto – 1994

- *QED is the very tentative title of a project to build a computer system that effectively represents all important mathematical knowledge and techniques.*
- *The QED system will conform to the highest standards of mathematical rigor, including the use of strict formality in the internal representation of knowledge and the use of mechanical methods to check proofs of the correctness of all entries in the system.*
- *The QED project will be a major scientific undertaking requiring the cooperation and effort of hundreds of deep mathematical minds, considerable ingenuity by many computer scientists, and broad support and leadership from research agencies.*
-
- *Never happened, but inspired a lot of development – “QED Singularity”*

Intros to ITP Systems and Formal Math

- Hales's talk at Bourbaki seminar:
<https://www.youtube.com/watch?v=wgfbt-X28XQ>
- Harrison's article on formalization:
<http://www.cl.cam.ac.uk/~jrh13/papers/cacm.pdf>
- Harrison, Urban, Wiedijk: History of Interactive Theorem Proving:
<http://www.cl.cam.ac.uk/~jrh13/papers/joerg.html>
- ITP, CPP, IJCAR, CADE conferences

Bird's Eye View of ITP Systems by T. Hales



HOL Light

HOL Light has an exquisite minimal design. It has the smallest kernel of any system. John Harrison is the sole



Mizar

Once the clear front-runner, it now shows signs of age. Do not expect to understand the inner workings of this system unless you have been



Coq

Coq is built of modular components on a foundation of dependent type theory. This system has grown one PhD thesis at a time.



Isabelle

Designed for use with multiple foundational architectures, Isabelle's early development featured classical constructions in set theory. However,



Metamath

Does this really work? Defying expectations, Metamath seems to function shockingly well for those who are happy to live without plumbing.



Lean

Lean is ambitious, and it will be massive. Do not be fooled by the name. "Construction area keep out" signs are prominently posted on the perimeter fencing.

Irrationality of $\sqrt{2}$ (informal text)

tiny proof from Hardy & Wright, collected by F. Wiedijk:

Theorem 43 (Pythagoras' theorem). $\sqrt{2}$ is irrational.

The traditional proof ascribed to Pythagoras runs as follows. If $\sqrt{2}$ is rational, then the equation

$$a^2 = 2b^2 \tag{4.3.1}$$

is soluble in integers a, b with $(a, b) = 1$. Hence a^2 is even, and therefore a is even. If $a = 2c$, then $4c^2 = 2b^2$, $2c^2 = b^2$, and b is also even, contrary to the hypothesis that $(a, b) = 1$. \square

Irrationality of $\sqrt{2}$ (Formal Proof Sketch)

exactly the same text in Mizar syntax:

```
theorem Th43: :: Pythagoras' theorem
  sqrt 2 is irrational
proof
  assume sqrt 2 is rational;
  consider a,b such that
4_3_1: a^2 = 2*b^2 and
  a,b are relative prime;
  a^2 is even;
  a is even;
  consider c such that a = 2*c;
  4*c^2 = 2*b^2;
  2*c^2 = b^2;
  b is even;
  thus contradiction;
end;
```

Irrationality of $\sqrt{2}$ in HOL Light

```
let Sqrt_2_Irrational = prove
  (~rational(sqrt(&2)))`,
  SIMP_TAC[rational; real_abs; Sqrt_Pos_Le; REAL_POS] THEN
  REWRITE_TAC[NOT_EXISTS_THM] THEN REPEAT GEN_TAC THEN
  DISCH_THEN(CONJUNCTS_THEN2 ASSUME_TAC MP_TAC) THEN
  SUBGOAL_THEN (~((&p / &q) pow 2 = sqrt(&2) pow 2))`
    (fun th -> MESON_TAC[th]) THEN
  SIMP_TAC[Sqrt_Pow_2; REAL_POS; REAL_POW_DIV] THEN
  ASM_SIMP_TAC[REAL_EQ_LDIV_EQ; REAL_OF_NUM_LT; REAL_POW_LT;
    ARITH_RULE `0 < q <=> ~(q = 0)`] THEN
  ASM_MESON_TAC[NSqrt_2; REAL_OF_NUM_POW; REAL_OF_NUM_MUL; REAL_OF_NUM_EQ]];
```

Irrationality of $\sqrt{2}$ in Isabelle/HOL

```
theorem sqrt2_not_rational:
  "sqrt (real 2)  $\notin$   $\mathbb{Q}$ "
proof
  assume "sqrt (real 2)  $\in$   $\mathbb{Q}$ "
  then obtain m n :: nat where
    n_nonzero: "n  $\neq$  0" and sqrt_rat: "|sqrt (real 2)| = real m / real n"
    and lowest_terms: "gcd m n = 1" ..
  from n_nonzero and sqrt_rat have "real m = |sqrt (real 2)| * real n" by simp
  then have "real (m2) = (sqrt (real 2))2 * real (n2)"
    by (auto simp add: power2_eq_square)
  also have "(sqrt (real 2))2 = real 2" by simp
  also have "... * real (m2) = real (2 * n2)" by simp
  finally have eq: "m2 = 2 * n2" ..
  hence "2 dvd m2" ..
  with two_is_prime have dvd_m: "2 dvd m" by (rule prime_dvd_power_two)
  then obtain k where "m = 2 * k" ..
  with eq have "2 * n2 = 22 * k2" by (auto simp add: power2_eq_square mult_ac)
  hence "n2 = 2 * k2" by simp
  hence "2 dvd n2" ..
  with two_is_prime have "2 dvd n" by (rule prime_dvd_power_two)
  with dvd_m have "2 dvd gcd m n" by (rule gcd_greatest)
  with lowest_terms have "2 dvd 1" by simp
  thus False by arith
qed
```

Irrationality of $\sqrt{2}$ in Coq

```
Theorem irrational_sqrt_2: irrational (sqrt 2%nat).
intros p q H H0; case H.
apply (main_thm (Zabs_nat p)).
replace (Div2.double (q * q)) with (2 * (q * q));
  [idtac | unfold Div2.double; ring].
case (eq_nat_dec (Zabs_nat p * Zabs_nat p) (2 * (q * q))); auto; intros H1.
case (not_nm_INR _ _ H1); (repeat rewrite mult_INR).
rewrite <- (sqrt_def (INR 2)); auto with real.
rewrite H0; auto with real.
assert (q <> 0%R :=> R); auto with real.
field; auto with real; case p; simpl; intros; ring.
Qed.
```

Irrationality of $\sqrt{2}$ in Otter

Problem

```
set(auto).
set(ur_res).
assign(max_distinct_vars, 1).
list(usable).
x = x.
m(1,x) = x. %identity
m(x,1) = x.
m(x,m(y,z)) = m(m(x,y),z). %assoc
m(x,y) = m(y,x). %comm
m(x,y) != m(x,z) | y = z. %cancel
-d(x,y) | m(x,f(x,y)) = y. %divides
m(x,z) != y | d(x,y).
-d(2,m(x,y)) | d(2,x) | d(2,y). %2 prime
m(a,a) = m(2,m(b,b)). % a/b=sqrt(2)
-d(x,a) | -d(x,b) | x = 1. % a/b lowest
2 != 1.
end_of_list.
```

Proof

```
1 [] m(x,y)!=m(x,z)|y=z.
2 [] -d(x,y)|m(x,f(x,y))=y.
3 [] m(x,y)!=z|d(x,z).
4 [] -d(2,m(x,y))|d(2,x)|d(2,y).
5 [] -d(x,a)|-d(x,b)|x=1.
6 [] 2!=1.
7 [factor,4.2.3] -d(2,m(x,x))|d(2,x).
13 [] m(x,m(y,z))=m(m(x,y),z).
14 [copy,13,flip.1] m(m(x,y),z)=m(x,m(y,
16 [] m(x,y)=m(y,x).
17 [] m(a,a)=m(2,m(b,b)).
18 [copy,17,flip.1] m(2,m(b,b))=m(a,a).
30 [hyper,18,3] d(2,m(a,a)).
39 [para_from,18.1.1,1.1.1] m(a,a)!=m(2,
42 [hyper,30,7] d(2,a).
46 [hyper,42,2] m(2,f(2,a))=a.
48 [ur,42,5,6] -d(2,b).
50 [ur,48,7] -d(2,m(b,b)).
59 [ur,50,3] m(2,x)!=m(b,b).
60 [copy,59,flip.1] m(b,b)!=m(2,x).
145 [para_from,46.1.1,14.1.1.1,flip.1] m
189 [ur,60,39] m(a,a)!=m(2,m(2,x)).
190 [copy,189,flip.1] m(2,m(2,x))!=m(a,a
1261 [para_into,145.1.1.2,16.1.1] m(2,m(
1272 [para_from,145.1.1,190.1.1.2] m(2,m
1273 [binary,1272.1,1261.1] $F.
```

Today: Computers Checking Large Math Proofs



Scientists Deliver Formal Proof of Famous Kepler Conjecture

Jun 16, 2017 by News Staff / Source

◀ Previous | Next ▶

Published in
Mathematics

Tagged as
Johannes Kepler
Kepler conjecture

**Follow
You Might Like**



Researchers Develop First-Ever 3D Numerical Model of Melting Snowflake



Researchers Develop Mathematical Model for How Innovations

An international team of mathematicians led by University of Pittsburgh Professor **Thomas Hales** has delivered a formal proof of the **Kepler conjecture**, a famous problem in discrete geometry. The team's **paper** is published in the journal *Forum of Mathematics, Pi*.



LATEST NEWS



SPHERE Captures Young Exoplanet Beta Pictoris b Orbiting around Its Star

Nov 13, 2018 | Astronomy



Mirace eatoni: Newly-Discovered Cretaceous Bird Lived Among Dinosaurs, Was Strong Flier

Nov 13, 2018 | Paleontology



Juno Takes Closer Look at Jupiter's Magnificent, Swirling Clouds

Nov 13, 2018 | Space Exploration



Physicists Solve Structure of Unusually Complex Form of Nitrogen

Nov 13, 2018 | Physical Chemistry



Natural Compound Protects Hypertensive Rats against Heart Disease

Nov 13, 2018 | Medicine



Inventive Orangutans Make Hook Tools to Retrieve Food

Nov 12, 2018 | Biology



Researchers Find 40,000-Year-Old Figurative Paintings in Bornean Cave

Nov 12, 2018 | Archaeology



Hubble Sees Lensing Galaxy Cluster,

cdn.sci-news.com/images/enlarge3/image_4960e-Kepler-Conjecture.jpg

Big Example: The Flyspeck project

- Kepler conjecture (1611): The most compact way of stacking balls of the same size in space is a pyramid.



$$V = \frac{\pi}{\sqrt{18}} \approx 74\%$$

- Formal proof finished in 2014
- 20000 lemmas in geometry, analysis, graph theory
- All of it at <https://code.google.com/p/flyspeck/>
- All of it **computer-understandable and verified** in HOL Light:
- `polyhedron s /\ c face_of s ==> polyhedron c`
- However, this took **20 – 30 person-years!**

Big Formalizations

- Kepler Conjecture (Hales et al, 2014, HOL Light, Isabelle)
- Feit-Thompson (odd-order) theorem
 - Two graduate books
 - Gonthier et al, 2012, Coq
- Compendium of Continuous Lattices (CCL)
 - 60% of the book formalized in Mizar
 - Bancerek, Trybulec et al, 2003
- The Four Color Theorem (Gonthier and Werner, 2005, Coq)

Mid-size Formalizations

- Gödel's First Incompleteness – N. Shankar (NQTHM), R. O'Connor (Coq)
- Brouwer Fixed Point Theorem – K. Pak (Mizar), J. Harrison (HOL Light)
- Jordan Curve Th. – T. Hales (HOL Light), A. Kornilowicz et al. (Mizar)
- Prime Number Th. – J. Avigad et al (Isab/HOL), J. Harrison (HOL Light)
- Gödel's Second Incompleteness Theorem – L. Paulson (Isabelle/HOL)
- Central Limit Theorem – J. Avigad (Isabelle/HOL)
- Consistency of the Negation of CH – J. Han and F. van Doorn (Lean)
- ... and many more

Large Software Verifications

- seL4 – operating system microkernel
 - Gerwin Klein and his group at NICTA, Isabelle/HOL
- CompCert – a formally verified C compiler
 - Xavier Leroy and his group at INRIA, Coq
- EURO-MILS – verified virtualization platform
 - ongoing 6M EUR FP7 project, Isabelle
- CakeML – verified implementation of ML
 - Magnus Myreen, HOL4

What Are Automated Theorem Provers?

- Computer programs that (try to) determine if
 - A conjecture C is a logical consequence of a set of axioms Ax
 - The derivation of conclusions that follow inevitably from facts.
- Systems: Vampire, E, SPASS, Prover9, Z3, CVC4, Satallax, iProver, ...
- Brute-force search calculi (resolution, superposition, tableaux, SMT, ...)
- Human-designed heuristics for pruning of the search space
- Combinatorial explosion on large KBs like Flyspeck and Mizar
- Need to be equipped with good domain-specific inference guidance ...
- ... and that is what we try to do ...
- ... typically by learning in various ways from over knowledge bases ...

Outline

Motivation, Learning vs. Reasoning

Computer Understandable (Formal) Math

Learning of Theorem Proving - Overview

Demos

High-level Reasoning Guidance: Premise Selection

Low Level Guidance of Theorem Provers

Mid-level Reasoning Guidance

Autoformalization

Using Learning to Guide Theorem Proving

- **high-level**: pre-select lemmas from a large library, give them to ATPs
- **high-level**: pre-select a good ATP strategy/portfolio for a problem
- **high-level**: pre-select good *hints* for a problem, use them to guide ATPs
- **low-level**: guide every inference step of ATPs (tableau, superposition)
- **low-level**: guide every kernel step of LCF-style ITPs
- **mid-level**: guide application of tactics in ITPs, learn new tactics
- **mid-level**: invent suitable strategies/procedures for classes of problems
- **mid-level**: invent suitable conjectures for a problem
- **mid-level**: invent suitable concepts/models for problems/theories
- **proof sketches**: explore stronger/related theories to get proof ideas
- **theory exploration**: develop interesting theories by conjecturing/proving
- **feedback loops**: (dis)prove, learn from it, (dis)prove more, learn more, ...
- **autoformalization**: (semi-)automate translation from \LaTeX to formal
- ...

Large Datasets

- Mizar / MML / MPTP – since 2003
- MPTP Challenge (2006), MPTP2078 (2011), Mizar40 (2013)
- Isabelle (and AFP) – since 2005
- Flyspeck (including core HOL Light and Multivariate) – since 2012
- HOL4 – since 2014, CakeML – 2017, GRUNGE – 2019
- Coq – since 2013/2016
- ACL2 – 2014?
- Lean?, Stacks?, Arxiv?, ProofWiki?, ...

Outline

Motivation, Learning vs. Reasoning

Computer Understandable (Formal) Math

Learning of Theorem Proving - Overview

Demos

High-level Reasoning Guidance: Premise Selection

Low Level Guidance of Theorem Provers

Mid-level Reasoning Guidance

Autoformalization

- **Hammering Mizar:** <http://grid01.ciirc.cvut.cz/~mptp/out4.ogv>
- **TacticToe on HOL4:**
http://grid01.ciirc.cvut.cz/~mptp/tactictoe_demo.ogv
- **Inf2formal over HOL Light:**
<http://grid01.ciirc.cvut.cz/~mptp/demo.ogv>
- **TacticToe longer:** <https://www.youtube.com/watch?v=B04Y8ynwT6Y>

Outline

Motivation, Learning vs. Reasoning

Computer Understandable (Formal) Math

Learning of Theorem Proving - Overview

Demos

High-level Reasoning Guidance: Premise Selection

Low Level Guidance of Theorem Provers

Mid-level Reasoning Guidance

Autoformalization

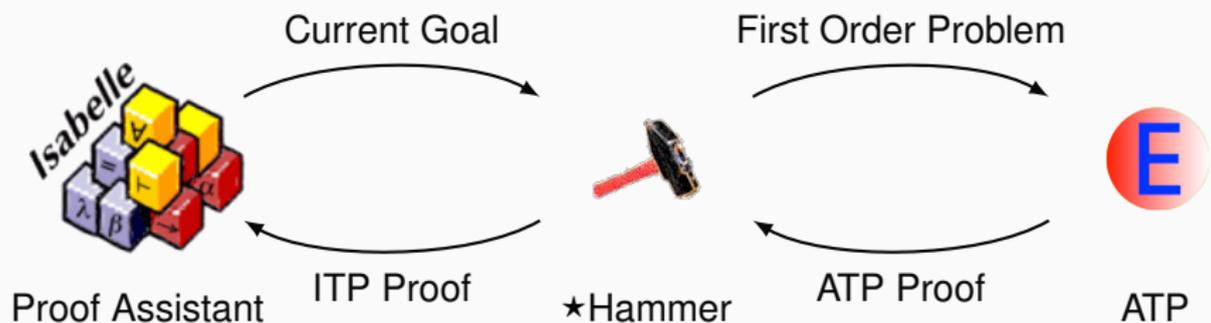
High-level ATP guidance: Premise Selection

- 2003: Can existing ATPs be used over the Mizar library?
- About 80000 nontrivial math facts at that time – impossible to use them all
- Is good premise selection for proving a new conjecture possible at all?
- Or is it a mysterious power of mathematicians? (Penrose)
- Today: Premise selection is not a mysterious property of mathematicians!
- Reasonably good algorithms started to appear (more below).

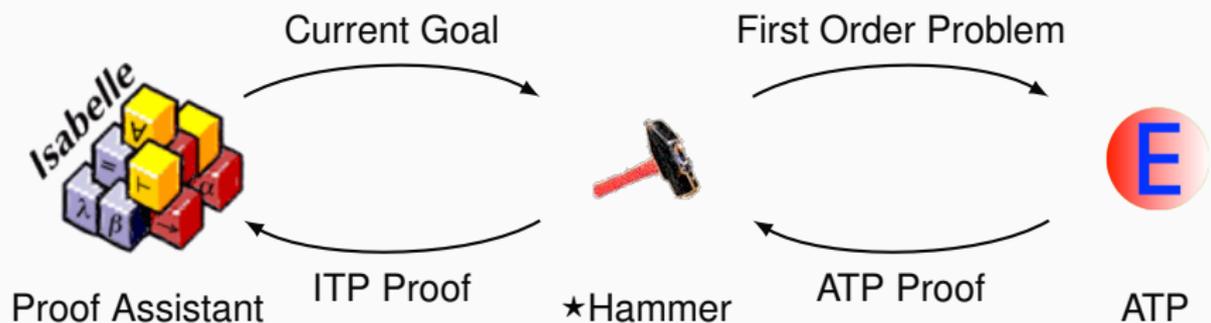
Example system: Mizar Proof Advisor (2003)

- train naive-Bayes fact selection on all previous Mizar/MML proofs (50k)
- input features: conjecture symbols; output labels: names of facts
- recommend relevant facts when proving new conjectures
- give them to unmodified FOL ATPs
- possibly reconstruct inside the ITP afterwards (lots of work)
- First results over the whole Mizar library in 2003:
 - about 70% coverage in the first 100 recommended premises
 - chain the recommendations with strong ATPs to get full proofs
 - about 14% of the Mizar theorems were then automatically provable (SPASS)
- Today's methods: about 45-50% (and we are still just beginning!)

Today's AI-ATP systems (★-Hammers)

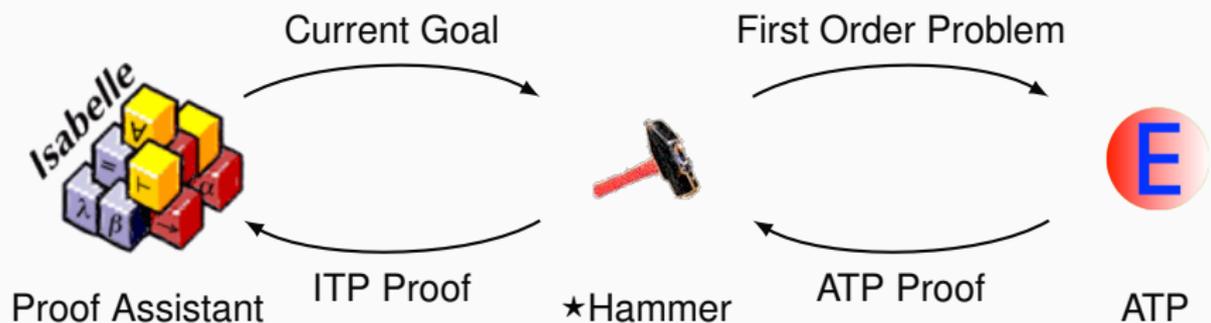


Today's AI-ATP systems (★-Hammers)



How much can it do?

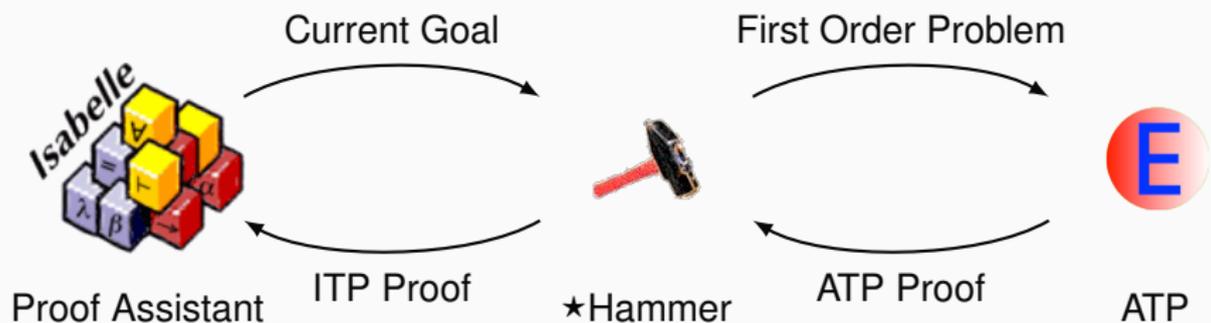
Today's AI-ATP systems (★-Hammers)



How much can it do?

- Mizar / MML – MizAR
- Isabelle (Auth, Jinja) – Sledgehammer
- Flyspeck (including core HOL Light and Multivariate) – HOL(y)Hammer
- HOL4 (Gauthier and Kaliszyk)
- CoqHammer (Czajka and Kaliszyk) - about 40% on Coq standard library

Today's AI-ATP systems (★-Hammers)



How much can it do?

- Mizar / MML – MizAR
- Isabelle (Auth, Jinja) – Sledgehammer
- Flyspeck (including core HOL Light and Multivariate) – HOL(y)Hammer
- HOL4 (Gauthier and Kaliszyk)
- CoqHammer (Czajka and Kaliszyk) - about 40% on Coq standard library

≈ 45% success rate

Various Improvements and Additions

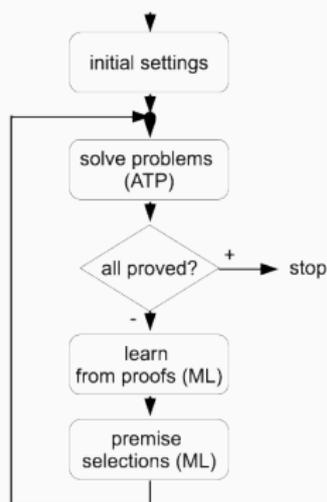
- Model-based features for **semantic similarity** [IJCAR'08]
- Features encoding **term matching/unification** [IJCAI'15]
- Stronger learners: SVMs, weighted k-NN, boosted trees (XGBoost)
- **Matching and transferring concepts** and theorems between libraries (Gauthier & Kaliszyk) – allows “superhammers”, conjecturing, and more
- Lemmatization – extracting and considering millions of low-level lemmas
- LSI, word2vec, neural models, definitional embeddings (with Google)
- Combined with tactical search and MCTS: **TacticToe** (Gauthier, 2017)
- Learning in binary setting from many alternative proofs
- Negative/positive mining (ATPBoost - Piotrowski & JU, 2018)

Summary of Features Used

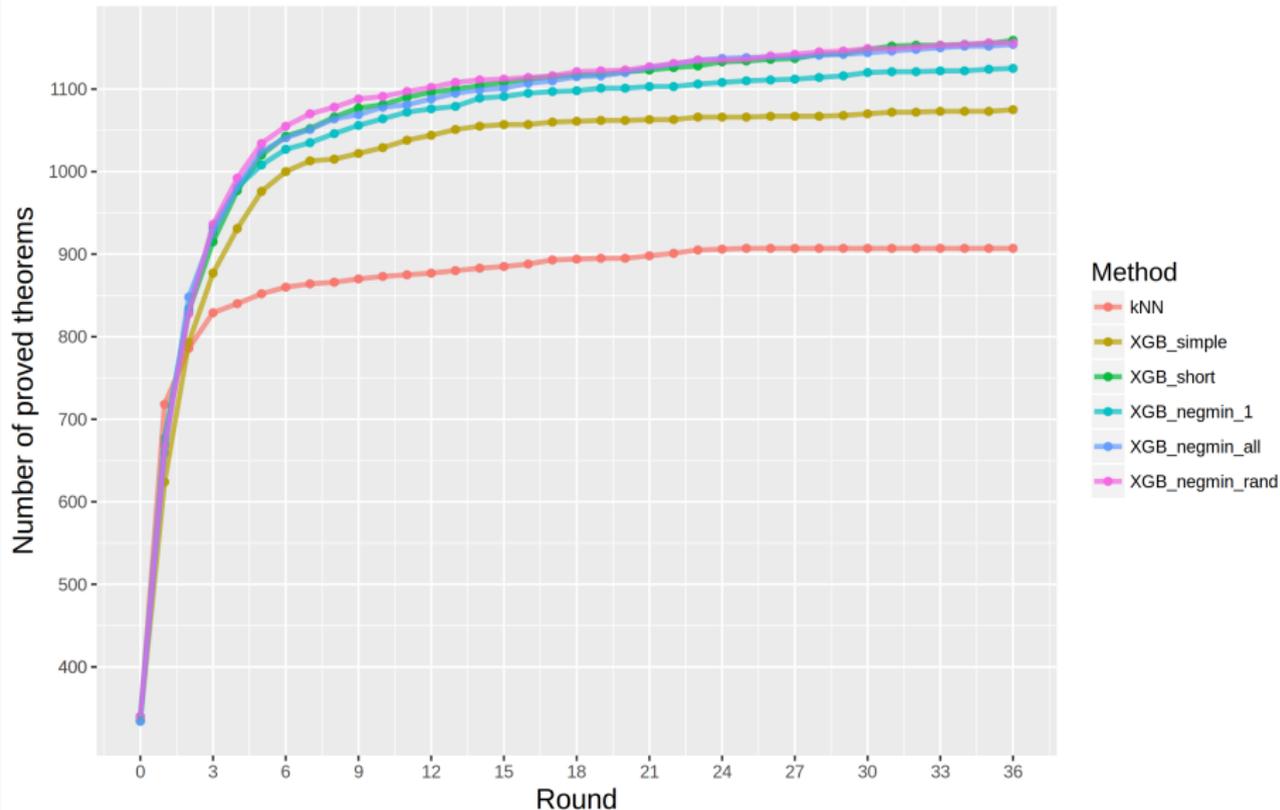
- From **syntactic** to more **semantic**:
- Constant and function symbols
- Walks in the term graph
- Walks in clauses with polarity and variables/skolems unified
- Subterms, de Bruijn normalized
- Subterms, all variables unified
- Matching terms, no generalizations
- terms and (some of) their generalizations
- Substitution tree nodes
- All unifying terms
- LSI/PCA, word2vec, fasttext, etc.
- Neural embeddings: CNN, RNN, Tree NN, Graph CNN, ...
- Evaluation in a large set of (finite) models
- Vectors of proof similarities (proof search hidden states)
- Vectors of problems solved (for ATP strategies)

High-level Feedback Loops – MALARea

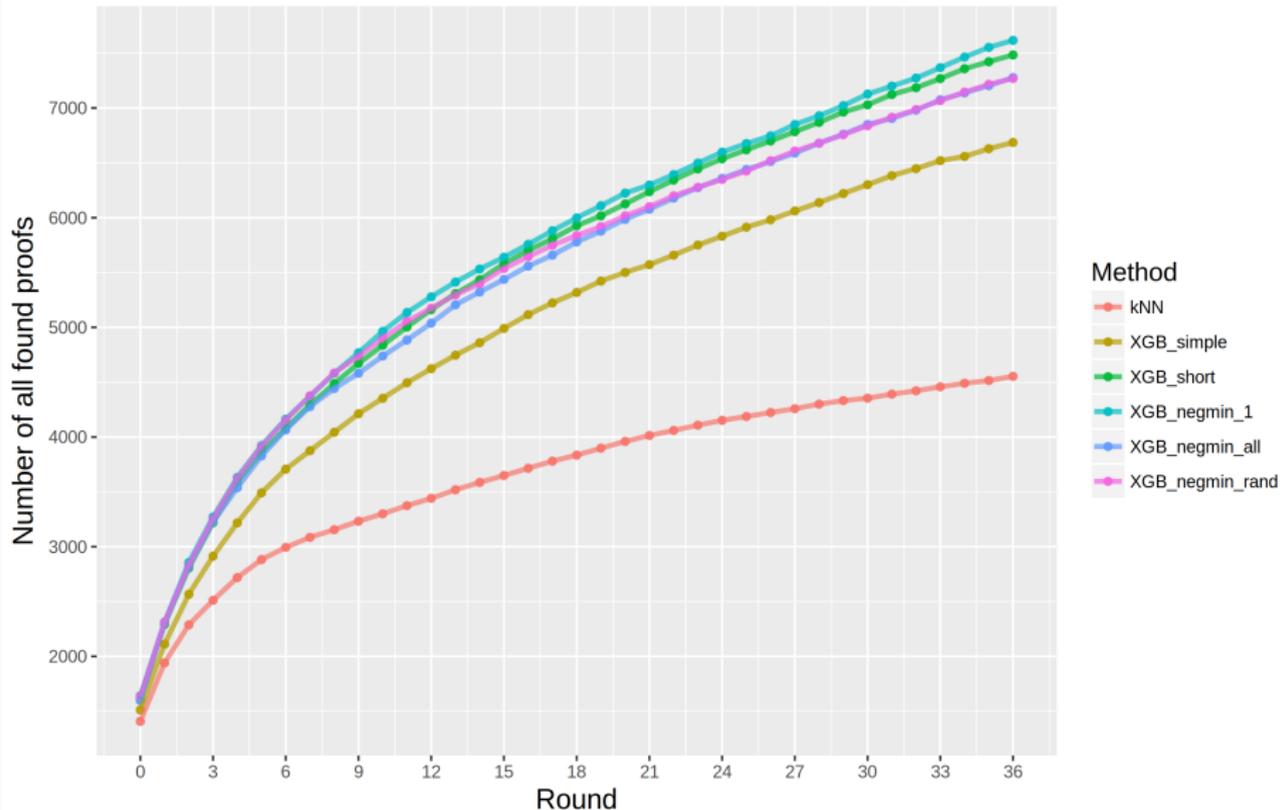
- Machine Learner for Autom. Reasoning (2006) – infinite hammering
- Feedback loop interleaving ATP with learning premise selection
- Both syntactic and **semantic** features for characterizing formulas:
- Evolving set of finite (counter)models in which formulas evaluated
- Winning AI/ATP benchmarks (MPTPChallenge, CASC 2008/12/13/18)
- ATPBoost (Piotrowski) - recent incarnation focusing on multiple proofs



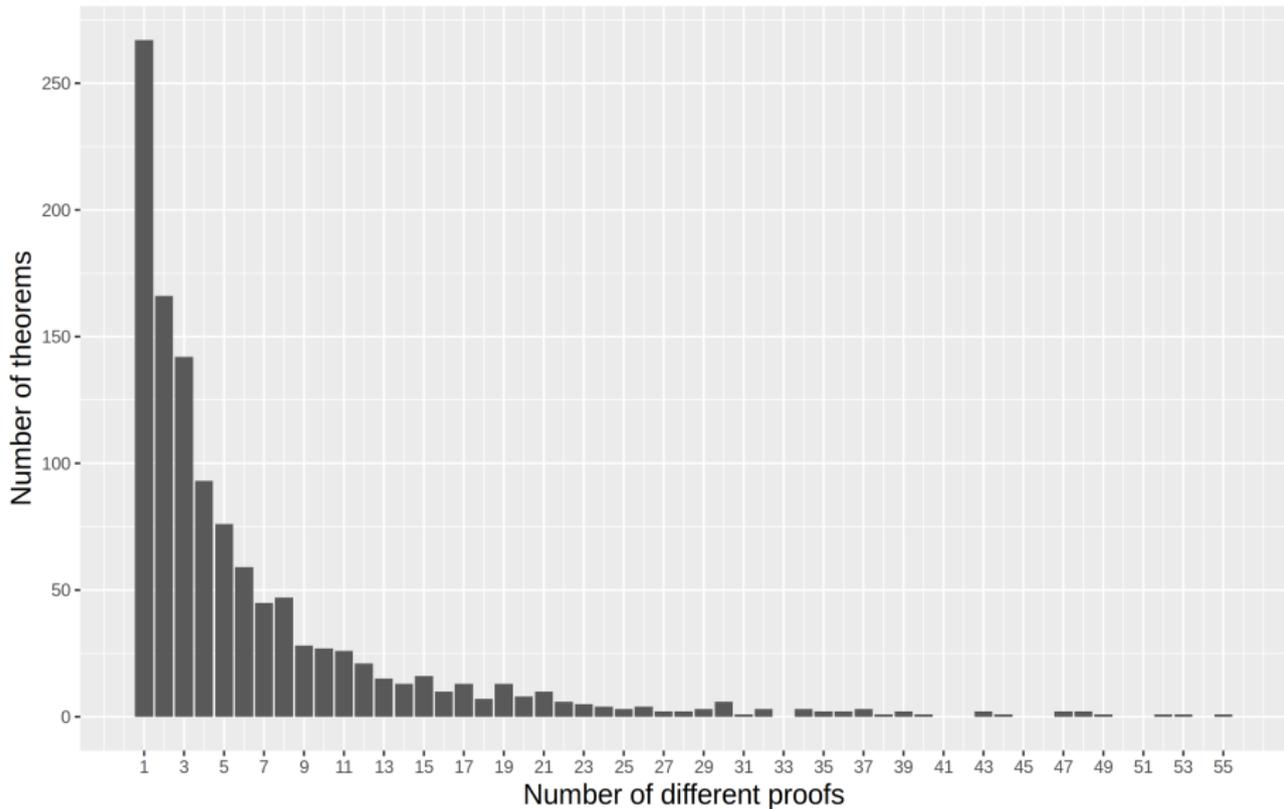
Prove-and-learn loop on MPTP2078 data set



Prove-and-learn loop on MPTP2078 data set



Number of found proofs per theorem at the end of the loop



Outline

Motivation, Learning vs. Reasoning

Computer Understandable (Formal) Math

Learning of Theorem Proving - Overview

Demos

High-level Reasoning Guidance: Premise Selection

Low Level Guidance of Theorem Provers

Mid-level Reasoning Guidance

Autoformalization

Low-level: Statistical Guidance of Connection Tableau

- learn guidance of every clausal inference in connection tableau (leanCoP)
- set of first-order clauses, *extension* and *reduction* steps
- proof finished when all branches are closed
- a lot of nondeterminism, requires backtracking
- *Iterative deepening* used in leanCoP to ensure completeness
- good for learning – the tableau compactly represents the proof state

Clauses:

$$c_1 : P(x)$$

$$c_2 : R(x, y) \vee \neg P(x) \vee Q(y)$$

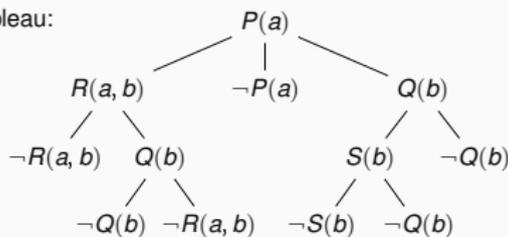
$$c_3 : S(x) \vee \neg Q(b)$$

$$c_4 : \neg S(x) \vee \neg Q(x)$$

$$c_5 : \neg Q(x) \vee \neg R(a, x)$$

$$c_6 : \neg R(a, x) \vee Q(x)$$

Closed Connection Tableau:



leanCoP: Minimal FOL Theorem Prover

```
% prove (Cla , Path , PathLim , Lem , Set)
prove ([ Lit | Cla ] , Path , PathLim , Lem , Set) :-
  ( - NegLit = Lit ; - Lit = NegLit ) ->
  (
    member (NegL , Path) ,
    unify_with_occurs_check (NegL , NegLit)
  ;
    % main nondeterminism
    lit (NegLit , NegL , Cla1 , Grnd1) ,
    unify_with_occurs_check (NegL , NegLit) ,
    prove (Cla1 , [ Lit | Path ] , PathLim , Lem , Set)
  ) ,
  prove (Cla , Path , PathLim , Lem , Set) .
prove ([ ] , _ , _ , _ , _) .
```

Statistical Guidance of Connection Tableau – rICoP

- **MaLeCoP** (2011): first prototype Machine Learning Connection Prover
- Fairly Efficient MaLeCoP = **FEMaLeCoP** (15% better than leanCoP)
- 2018: remove iterative deepening, the prover can go deep (completeness bad!)
- Monte-Carlo Tree Search (MCTS) governs the search (AlphaGo/Zero!)
- MCTS search nodes are sequences of clause application
- a good heuristic to explore new vs exploit good nodes:

$$\frac{w_i}{n_i} + c \cdot p_i \cdot \sqrt{\frac{\ln N}{n_i}} \quad (\text{UCT - Kocsis, Szepesvari 2006})$$

- learning both *policy* (clause selection) and *value* (state evaluation)
- clauses represented not by names but also by features (generalize!)
- **binary** learning setting used: | proof state | clause features |
- mostly term walks of length 3 (trigrams), hashed into small integers
- **fast strong learners** via C interface to OCAML (boosted trees)
- many iterations of proving and learning

Statistical Guidance of Connection Tableau – rICoP

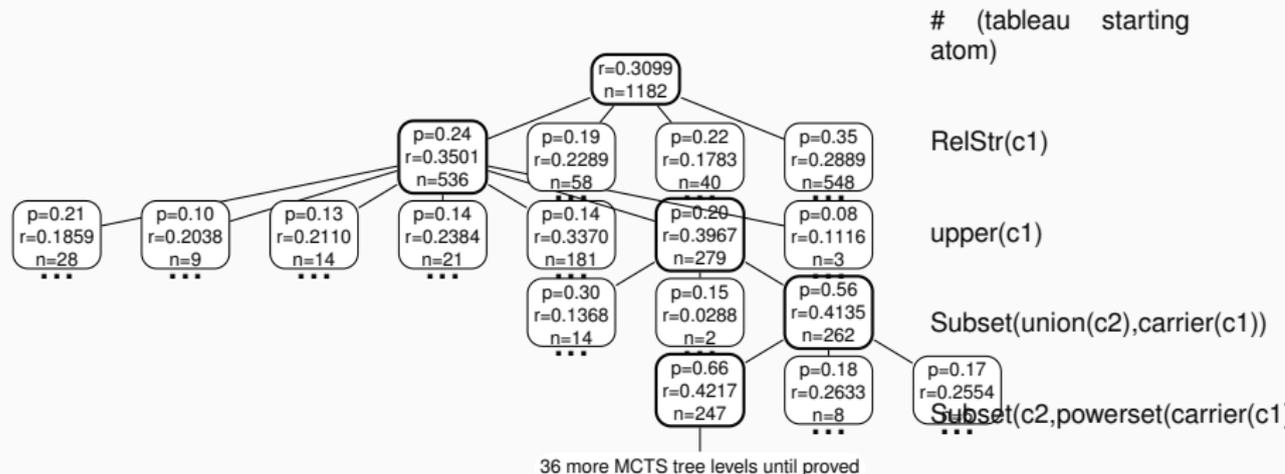
- On 32k Mizar40 problems using 200k inference limit
- nonlearning CoPs:

System	leanCoP	bare prover	rICoP no policy/value (UCT only)
Training problems proved	10438	4184	7348
Testing problems proved	1143	431	804
Total problems proved	11581	4615	8152

- rICoP with policy/value after 5 proving/learning iters on the training data
- $1624/1143 = 42.1\%$ improvement over leanCoP on the testing problems

Iteration	1	2	3	4	5	6	7	8
Training proved	12325	13749	14155	14363	14403	14431	14342	14498
Testing proved	1354	1519	1566	1595	1624	1586	1582	1591

Tree Example



Recent Variations – FLoP, RNN, GNN

- Finding Longer Proofs – Zombori & al <https://arxiv.org/abs/1905.13100>
- Curriculum Learning in connection tableau over Robinson Arithmetic
- Addition and multiplication learned perfectly from $1 * 1 = 1$
- Headed towards learning algorithms/decision procedures from math data
- Now black-box. Combining with symbolic methods (ILP?) our next target
- Using RNNs for better tableau encoding, prediction of actions ...
- ... even guessing (decoding) next tableau literals (Piotrowski & JU, 2019)
- The same with graph neural nets (GNN) - encouraging preliminary results

Side Note on Symbolic Learning with NNs

- Recurrent NNs with attention recently very good at the inf2formal task
- Experiments with using them for **symbolic rewriting** (Piotrowski et. al)
- We can **learn rewrite rules** from sufficiently many data
- 80-90% on algebra datasets, 70-99% on normalizing polynomials
- complements symbolic methods like ILP that suffer if too much data
- Similar use for **conjecturing** (Chvalovsky et al):
- Learn *consistent translations* between different math contexts:
- additive groups \rightarrow multiplicative groups

Side Note on Symbolic Learning with NNs

Table: Examples in the AIM data set.

Rewrite rule:	Before rewriting:	After rewriting:
$b(s(e, v1), e) = v1$	$k(b(s(e, v1), e), v0)$	$k(v1, v0)$
$o(v0, e) = v0$	$t(v0, o(v1, o(v2, e)))$	$t(v0, o(v1, v2))$

Table: Examples in the polynomial data set.

Before rewriting:	After rewriting:
$(x * (x + 1)) + 1$	$x^2 + x + 1$
$(2 * y) + 1 + (y * y)$	$y^2 + 2 * y + 1$
$(x + 2) * ((2 * x) + 1) + (y + 1)$	$2 * x^2 + 5 * x + y + 3$

Side Note on Conjecturing with RNNs

- Use HO abstraction to align concepts with similar properties
- Train RNNs to learn the translations between different contexts

Example: obtain a new valid automatically provable lemma

$$(X \cap Y) \setminus Z = (X \setminus Z) \cap (Y \setminus Z)$$

from

$$(X \cup Y) \setminus Z = (X \setminus Z) \cup (Y \setminus Z)$$

Examples of false but syntactically consistent conjectures:

for n, m being natural numbers holds $n \text{ gcd } m = n \text{ div } m$;

for R being Relation holds

$\text{with_suprema}(A) \iff \text{with_suprema}(\text{inverse_relation}(A))$;

Side Note on Model Learning with NNs

- Smolik 2019 (MSc thesis): modelling mathematical structures with NNs
- NNs reasonably learn cyclic groups and their extensions
- ... so far struggle in learning bigger permutation groups
- Plan: learn composition/variation of complicated math structures
- Use for model-style evaluation of formulas, conjectures, etc. – similarly to the finite models in Malarea, etc.

Guiding Saturation-Based Theorem Proving

Basic Saturation Loop – Given Clause Loop

```
 $P := \emptyset$  (processed)  
 $U := \{\textit{classified axioms and a negated conjecture}\}$  (unprocessed)  
while ( $U \neq \emptyset$ ) do  
  if ( $\perp \in U \cup P$ ) then return Unsatisfiable  
   $g := \textit{select}(U)$  (choose a given clause)  
   $P := P \cup \{g\}$  (add to processed)  
   $U := U \setminus \{g\}$  (remove from unprocessed)  
   $U := U \cup \{\textit{all clauses inferred from } g \textit{ and } P\}$  (add inferences)  
done  
return Satisfiable
```

Typically, U grows quadratically wrt. P

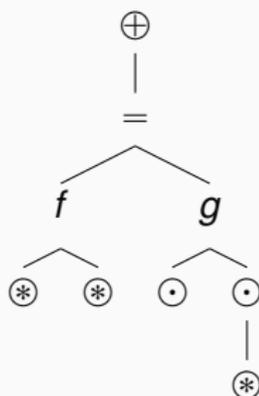
1M clauses in U in 10s common – choosing good g gets harder

Statistical Guidance of the Given Clause in E

- harder for learning than tableau, but **state-of-the-art in ATP**
- the proof state are two **large heaps of clauses** *processed/unprocessed*
- 2017: ENIGMA - manual feature engineering (Jakubuv & JU 2017)
- 2017: Deep guidance (neural nets) (Loos et al. 2017)
- both learn on E's proof search traces, put classifier in E
- positive examples: given clauses used in the proof
- negative examples: given clauses not used in the proof
- ENIGMA: fast feature extraction followed by fast/sparse linear classifier
- about 80% improvement on the AIM benchmark
- Deep guidance: convolutional nets - no feature engineering but slow
- ENIGMA-NG: better features and ML, gradient-boosted trees, tree NNs
- NNs made competitive in real-time, boosted trees still best

Clauses as Feature Vectors for ENIGMA

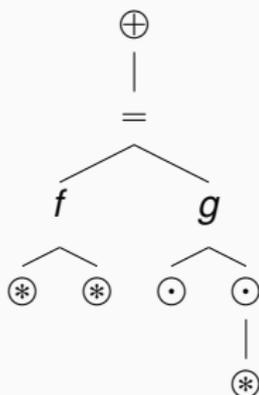
Collect and enumerate all the features. Count the clause features.



#	feature	count
1	($\oplus, =, a$)	0
\vdots	\vdots	\vdots
11	($\oplus, =, f$)	1
12	($\oplus, =, g$)	1
13	($=, f, *$)	2
14	($=, g, \odot$)	2
15	($g, \odot, *$)	1
\vdots	\vdots	\vdots

Clauses as Feature Vectors for ENIGMA

Take the counts as a **feature vector**.



#	feature	count
1	($\oplus, =, a$)	0
⋮	⋮	⋮
11	($\oplus, =, f$)	1
12	($\oplus, =, g$)	1
13	($=, f, *$)	2
14	($=, g, \odot$)	2
15	($g, \odot, *$)	1
⋮	⋮	⋮

Feedback loop for ENIGMA on Mizar data

- Similar to rICoP - interleave proving and learning of ENIGMA guidance
- Done on 57880 Mizar problems very recently
- **Ultimately a 70% improvement over the original strategy**
- Example Mizar proof found by ENIGMA: http://grid01.ciirc.cvut.cz/~mptp/7.13.01_4.181.1147/html/knaster#T21
- Its E-ENIGMA proof: http://grid01.ciirc.cvut.cz/~mptp/t21_knaster

	S	$S \odot M_9^0$	$S \oplus M_9^0$	$S \odot M_9^1$	$S \oplus M_9^1$	$S \odot M_9^2$	$S \oplus M_9^2$	$S \odot M_9^3$	$S \oplus M_9^3$
solved	14933	16574	20366	21564	22839	22413	23467	22910	23753
$S\%$	+0%	+10.5%	+35.8%	+43.8%	+52.3%	+49.4%	+56.5%	+52.8%	+58.4
$S+$	+0	+4364	+6215	+7774	+8414	+8407	+8964	+8822	+9274
$S-$	-0	-2723	-782	-1143	-508	-927	-430	-845	-454

	$S \odot M_{12}^3$	$S \oplus M_{12}^3$	$S \odot M_{16}^3$	$S \oplus M_{16}^3$
solved	24159	24701	25100	25397
$S\%$	+61.1%	+64.8%	+68.0%	+70.0%
$S+$	+9761	+10063	+10476	+10647
$S-$	-535	-295	-309	-183

ENIGMA Proof Example – Knaster

```
theorem Th21:
  ex a st a is_a_fixpoint_of f
proof
  set H = {h where h is Element of L: h [= f.h];
  set fH = {f.h where h is Element of L: h [= f.h];
  set uH = "\/"(H, L);
  set fuH = "\/"(fH, L);
  take uH;
  now
    let fh be Element of L;
    assume fh in fH;
    then consider h being Element of L such that
A1: fh = f.h and
A2: h [= f.h;
    h in H by A2;
    then h [= uH by LATTICE3:38;
    hence fh [= f.uH by A1,QUANTAL1:def 12;
  end;
  then fH is_less_than f.uH by LATTICE3:def 17;
  then
A3: fuH [= f.uH by LATTICE3:def 21;
  now
    let a be Element of L;
    assume a in H;
    then consider h being Element of L such that
A4: a = h & h [= f.h;
    reconsider fh = f.h as Element of L;
    take fh;
    thus a [= fh & fh in fH by A4;
  end;
  then uH [= fuH by LATTICE3:47;
  then
A5: uH [= f.uH by A3,LATTICES:7;
  then f.uH [= f.(f.uH) by QUANTAL1:def 12;
  then f.uH in H;
  then f.uH [= uH by LATTICE3:38;
  hence uH = f.uH by A5,LATTICES:8;
end;
```

ENIGMA Machine Learning Behavior

	<i>Liblinear</i>	<i>Xgboost</i>	<i>TreeNN</i>
TPR	90.54 %	99.36 %	97.82 %
TNR	83.52 %	93.32 %	94.69 %

Table: True Positive Rate (TPR) and True Negative Rate (TNR) on training data.

	<i>Liblinear</i>	<i>Xgboost</i>	<i>TreeNN</i>
TPR	80.54 %	83.35 %	82.00 %
TNR	62.28 %	72.60 %	76.88 %

Table: TPR and TNR on newly solved problems.

- TreeNNs about 10 times slower than XGBoost. XGBoost has 60% of the speed of unmodified E. Liblinear 90%.
- Fast feature hashing: immediate dimensionality reduction to 32k - 2k buckets. Good with XGBoost.

ProofWatch: Symbolic/Statistical Guidance of E

- Bob Veroff's *hints method* used for Prover9
- solve many easier/related problems, produce millions of lemmas
- load the useful lemmas on the *watchlist* (kind of conjecturing)
- *boost inferences on clauses that subsume a watchlist* clause
- watchlist parts are *fast thinking*, bridged by *standard (slow) search*
- *symbolic guidance*, initial attempts to choose good hints by statistical ML
- Very *long proofs of open conjectures* in quasigroup/loop theory (AIM)
- ProofWatch (Goertzel et al. 2018): load many proofs separately
- *dynamically* boost those that have been covered more
- needed for heterogeneous ITP libraries
- *statistical*: watchlists chosen using similarity and usefulness
- *semantic/deductive*: dynamic guidance based on exact proof matching
- results in *better vectorial characterization* of saturation proof searches

ProofWatch: Statistical/Symbolic Guidance of E

```
theorem Th36: :: YELLOW_5:36
```

```
for L being non empty Boolean RelStr for a, b being Element of L  
holds ( 'not' (a "\/" b) = ('not' a) "\/" ('not' b)  
      & 'not' (a "\/" b) = ('not' a) "\/" ('not' b) )
```

- De Morgan's laws for Boolean lattices
- guided by 32 related proofs resulting in 2220 watchlist clauses
- 5218 given clause loops, resulting ATP proof is 436 clauses
- 194 given clauses match the watchlist and 120 (61.8%) used in the proof
- most helped by the proof of WAYBEL_1:85 – done for lower-bounded Heyting

```
theorem :: WAYBEL_1:85
```

```
for H being non empty lower-bounded RelStr st H is Heyting holds  
for a, b being Element of H holds  
'not' (a "\/" b) >= ('not' a) "\/" ('not' b)
```

ProofWatch: Vectorial Proof State

Final state of the proof progress for the 32 proofs guiding YELLOW_5 : 36

0	0.438	42/96	1	0.727	56/77	2	0.865	45/52	3	0.360	9/25
4	0.750	51/68	5	0.259	7/27	6	0.805	62/77	7	0.302	73/242
8	0.652	15/23	9	0.286	8/28	10	0.259	7/27	11	0.338	24/71
12	0.680	17/25	13	0.509	27/53	14	0.357	10/28	15	0.568	25/44
16	0.703	52/74	17	0.029	8/272	18	0.379	33/87	19	0.424	14/33
20	0.471	16/34	21	0.323	20/62	22	0.333	7/21	23	0.520	26/50
24	0.524	22/42	25	0.523	45/86	26	0.462	6/13	27	0.370	20/54
28	0.411	30/73	29	0.364	20/55	30	0.571	16/28	31	0.357	10/28

EnigmaWatch: ProofWatch used with ENIGMA

- Use the **proof completion ratios as features** for **characterizing proof state**
- Instead of just static conjecture features - **the proof vectors evolve**
- Feed them to ML systems along with other features
- Good improvement, extendable in various ways
- Alternative to backtracking-based RL-style systems? (**Doesn't forget!**)

EnigmaWatch: ProofWatch used with ENIGMA

Baseline	Mean	Var	Corr	Rand	Baseline \cup Mean	Total
1140	1357	1345	1337	1352	1416	1483

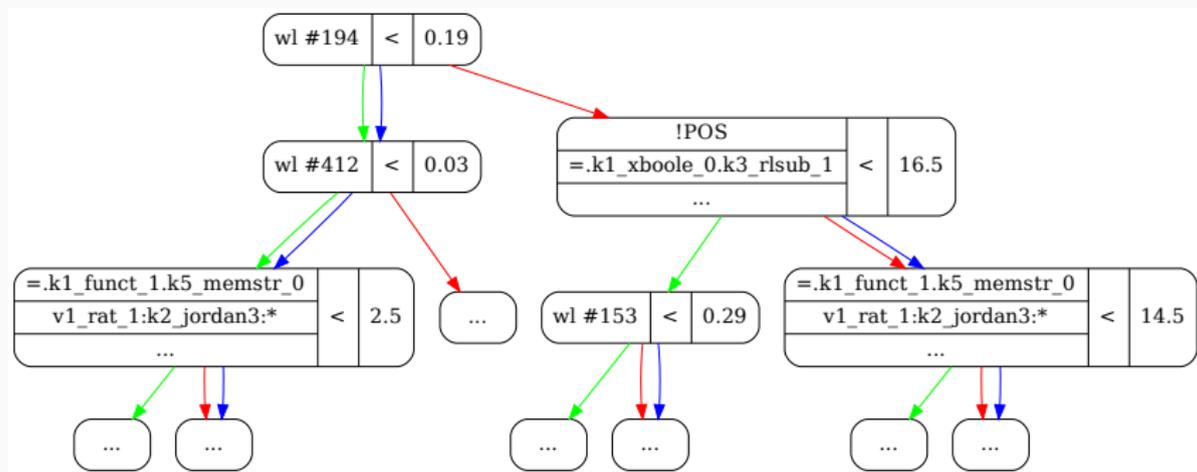
Table: ProofWatch evaluation: Problems solved by different versions.

loop	ENIGMA	Mean	Var	Corr	Rand	ENIGMA \cup Mean	Total
0	1557	1694	1674	1665	1690	1830	1974
1	1776	1815	1812	1812	1847	1983	2131
2	1871	1902	1912	1882	1915	2058	2200
3	1931	1954	1946	1920	1926	2110	2227

Table: ENIGMAWatch evaluation: Problems solved and the effect of looping.

- ENIGMAWatch initially much better, ENIGMA eventually catches up
- Still, EW produces simpler XGBoost models (easier learning)

Example of an XGBoost decision tree



Outline

Motivation, Learning vs. Reasoning

Computer Understandable (Formal) Math

Learning of Theorem Proving - Overview

Demos

High-level Reasoning Guidance: Premise Selection

Low Level Guidance of Theorem Provers

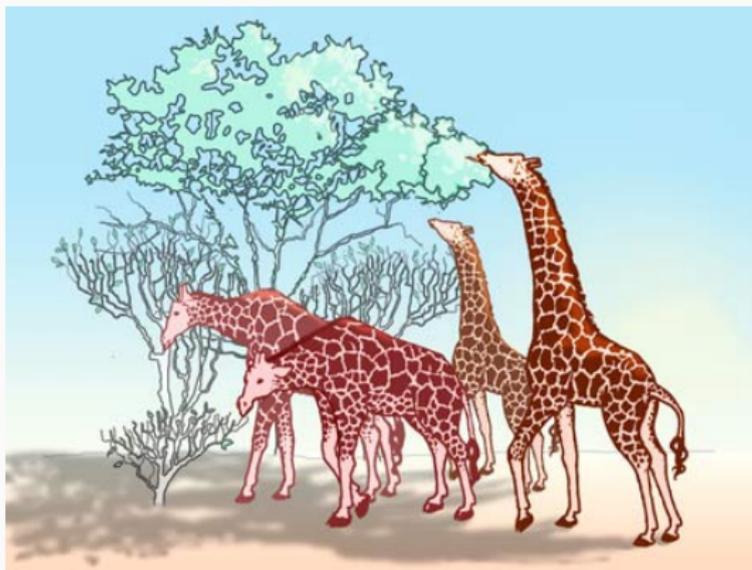
Mid-level Reasoning Guidance

Autoformalization

TacticToe: mid-level ITP Guidance (Gauthier et al.'18)

- learns from human **tactical HOL4 proofs** to solve new goals
- also from its own proofs - feedback loops
- no translation or reconstruction needed
- similar to rlCoP: policy/value learning
- however **much more technically challenging**:
 - tactic and goal state recording
 - tactic argument abstraction
 - absolutization of tactic names
 - nontrivial evaluation issues
 - learning of new tactics
- **policy**: which tactic/parameters to choose for a current goal?
- **value**: how likely is this proof state succeed?
- **66%** of HOL4 toplevel proofs in 60s (better than a hammer!)
- similar recent work for **Isabelle** (Nagashima 2018)
- work in progress for **Coq** and **HOL Light** (several groups)
- A lot of possible extensions ...

More Mid-level guidance: BliStr: Blind Strategymaker



- ATP **strategies are programs** specified in rich DSLs - can be **evolved**
- The ATP strategies are like giraffes, the problems are their food
- The better the giraffe specializes for eating problems unsolvable by others, the more it gets fed and further evolved

The E strategy with longest specification in Jan 2012

G-E--_029_K18_F1_PI_AE_SU_R4_CS_SP_S0Y:

```
4 * ConjectureGeneralSymbolWeight (
    SimulateSOS,100,100,100,50,50,10,50,1.5,1.5,1),
3 * ConjectureGeneralSymbolWeight (
    PreferNonGoals,200,100,200,50,50,1,100,1.5,1.5,1),
1 * Clauseweight (PreferProcessed,1,1,1),
1 * FIFOWeight (PreferProcessed)
```

BliStr: Blind Strategymaker

- **Strategies** characterized by the problems they solve
- **Problems** characterized by the strategies that solve them
- Improve on sets of **similar easy** problems to train for **unsolved** problems
- Interleave **low-time training on easy problems** with **high-time evaluation**
- Single strategy evolution done by **ParamILS** - Iterated Local Search (Hutter et al. 2009 – genetic methods work too)
- Thus **co-evolve** the strategies and their training problems
- The hard problems gradually become easier and turn into training data (the trees get lower for a taller giraffe)
- In the end, learn which strategy to use on which problem

The Longest E Strategy After Evolution

atpstr_my_c7bb78cc4c665670e6b866a847165cb4bf997f8a:

```
6 * ConjectureGeneralSymbolWeight (PreferNonGoals,100,100,100,50,50,1000,100,1.5,1.5,1)
8 * ConjectureGeneralSymbolWeight (PreferNonGoals,200,100,200,50,50,1,100,1.5,1.5,1)
8 * ConjectureGeneralSymbolWeight (SimulateSOS,100,100,100,50,50,50,50,1.5,1.5,1)
4 * ConjectureRelativeSymbolWeight (ConstPrio,0.1, 100, 100, 100, 100, 1.5, 1.5, 1.5)
10 * ConjectureRelativeSymbolWeight (PreferNonGoals,0.5, 100, 100, 100, 100, 1.5, 1.5, 1)
2 * ConjectureRelativeSymbolWeight (SimulateSOS,0.5, 100, 100, 100, 100, 1.5, 1.5, 1)
10 * ConjectureSymbolWeight (ConstPrio,10,10,5,5,5,1.5,1.5,1.5)
1 * Clauseweight (ByCreationDate,2,1,0.8)
1 * Clauseweight (ConstPrio,3,1,1)
6 * Clauseweight (ConstPrio,1,1,1)
2 * Clauseweight (PreferProcessed,1,1,1)
6 * FIFOWeight (ByNegLitDist)
1 * FIFOWeight (ConstPrio)
2 * FIFOWeight (SimulateSOS)
8 * OrientLMaxWeight (ConstPrio,2,1,2,1,1)
2 * PNRefinedweight (PreferGoals,1,1,1,2,2,2,0.5)
10 * RelevanceLevelWeight (ConstPrio,2,2,0,2,100,100,100,100,1.5,1.5,1)
8 * RelevanceLevelWeight2 (PreferNonGoals,0,2,1,2,100,100,100,400,1.5,1.5,1)
2 * RelevanceLevelWeight2 (PreferGoals,1,2,1,2,100,100,100,400,1.5,1.5,1)
6 * RelevanceLevelWeight2 (SimulateSOS,0,2,1,2,100,100,100,400,1.5,1.5,1)
8 * RelevanceLevelWeight2 (SimulateSOS,1,2,0,2,100,100,100,400,1.5,1.5,1)
5 * rweight21_g
3 * Refinedweight (PreferNonGoals,1,1,2,1.5,1.5)
1 * Refinedweight (PreferNonGoals,2,1,2,2,2)
2 * Refinedweight (PreferNonGoals,2,1,2,3,0.8)
8 * Refinedweight (PreferGoals,1,2,2,1,0.8)
10 * Refinedweight (PreferGroundGoals,2,1,2,1.0,1)
20 * Refinedweight (SimulateSOS,1,1,2,1.5,2)
1 * Refinedweight (SimulateSOS,3,2,2,1.5,2)
```

BliStr Evaluation on 1000 Mizar problems

- Original E coverage: 597 problems
- After 30 hours of strategy growing: 22 strategies covering 670 problems
- Best strategy solves 598 problems (1 more than all original strategies)
- Portfolio of 14 strategies improves E auto-mode by 25% on new problems
- Similar results for the Flyspeck problems
- Future: enrich the DSLs with new methods (ENIGMA/Watch, etc.)

Outline

Motivation, Learning vs. Reasoning

Computer Understandable (Formal) Math

Learning of Theorem Proving - Overview

Demos

High-level Reasoning Guidance: Premise Selection

Low Level Guidance of Theorem Provers

Mid-level Reasoning Guidance

Autoformalization

Autoformalization

- Goal: Learn **understanding of informal math** formulas and reasoning
- Experiments with the **CYK chart parser linked to semantic methods**
- Experiments with **neural methods**
- Semantic methods: Type checking, theorem proving
- Corpora: Flyspeck, Mizar, Proofwiki, etc.

Statistical/Semantic Parsing of Informalized HOL

- Training and testing examples exported from Flyspeck formulas
 - Along with their **informalized** versions
- Grammar parse trees
 - Annotate each (nonterminal) symbol with its **HOL type**
 - Also “semantic (formal)” nonterminals annotate overloaded terminals
 - guiding analogy: word-sense disambiguation using CYK is common
- Terminals exactly compose the textual form, for example:

- **REAL_NEGNEG**: $\forall x. --x = x$

```
(Comb (Const "!" (Tyapp "fun" (Tyapp "fun" (Tyapp "real") (Tyapp "bool"))  
(Tyapp "bool")))) (Abs "A0" (Tyapp "real") (Comb (Comb (Const "=" (Tyapp "fun"  
(Tyapp "real") (Tyapp "fun" (Tyapp "real") (Tyapp "bool")))) (Comb (Const  
"real_neg" (Tyapp "fun" (Tyapp "real") (Tyapp "real")))) (Comb (Const  
"real_neg" (Tyapp "fun" (Tyapp "real") (Tyapp "real")))) (Var "A0" (Tyapp  
"real")))) (Var "A0" (Tyapp "real"))))
```

- **becomes**

```
("(Type bool)" ! ("(Type (fun real bool))" (Abs ("(Type real)"  
(Var A0)) ("(Type bool)" ("(Type real)" real_neg ("(Type real)"  
real_neg ("(Type real)" (Var A0)))) = ("(Type real)" (Var A0))))))
```


CYK Learning and Parsing (KUV, ITP 17)

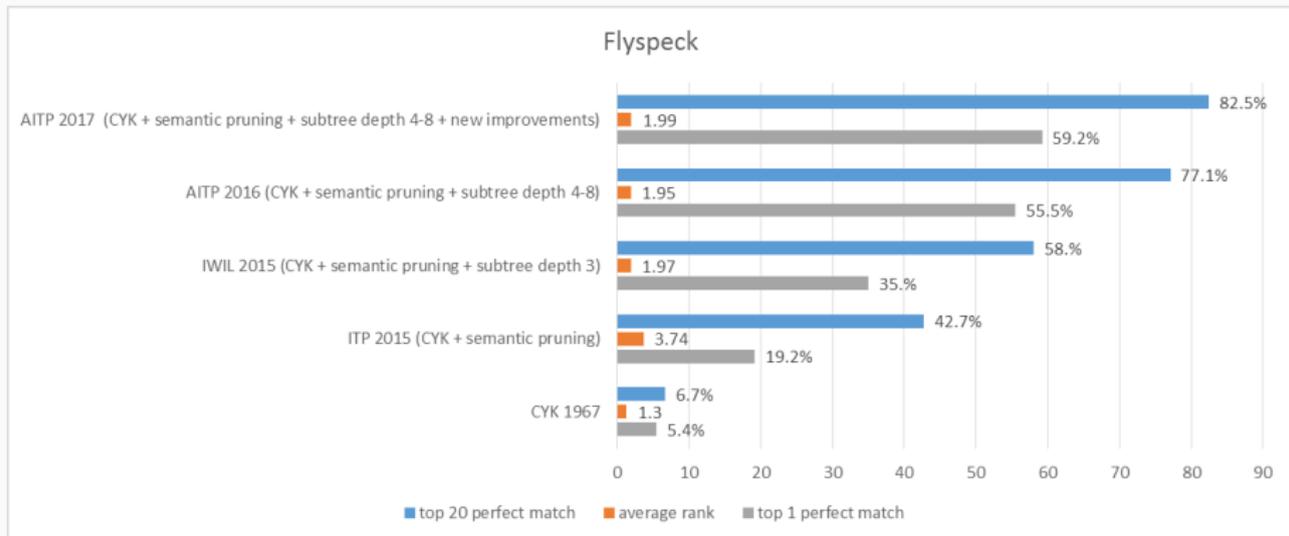
- Induce **PCFG** (probabilistic context-free grammar) from the trees
 - Grammar rules obtained from the inner nodes of each grammar tree
 - Probabilities are computed from the **frequencies**
- The PCFG grammar is binarized for efficiency
 - New nonterminals as shortcuts for multiple nonterminals
- CYK: dynamic-programming algorithm for parsing **ambiguous sentences**
 - input: sentence – a sequence of words and a binarized PCFG
 - output: N **most probable** parse trees
- Additional **semantic** pruning
 - **Compatible types for free variables in subtrees**
 - Optionally merge equivalent subtrees (tricky)
- Allow small probability for each symbol to be a variable
- Top parse trees are de-binarized to the original CFG
 - Transformed to HOL parse trees (preterms, Hindley-Milner)
 - typed checked in HOL and then given to an ATP (hammer)

Autoformalization based on PCFG and semantics

- “`sin (0 * x) = cos pi / 2`”
- produces 16 parses
- of which 11 get type-checked by HOL Light as follows
- with all but three being proved by HOL(y)Hammer
- **demo:** <http://grid01.ciirc.cvut.cz/~mptp/demo.ogv>

```
sin (&0 * A0) = cos (pi / &2) where A0:real
sin (&0 * A0) = cos pi / &2 where A0:real
sin (&0 * &A0) = cos (pi / &2) where A0:num
sin (&0 * &A0) = cos pi / &2 where A0:num
sin (&(0 * A0)) = cos (pi / &2) where A0:num
sin (&(0 * A0)) = cos pi / &2 where A0:num
csin (Cx (&0 * A0)) = ccos (Cx (pi / &2)) where A0:real
csin (Cx (&0) * A0) = ccos (Cx (pi / &2)) where A0:real^2
Cx (sin (&0 * A0)) = ccos (Cx (pi / &2)) where A0:real
csin (Cx (&0 * A0)) = Cx (cos (pi / &2)) where A0:real
csin (Cx (&0) * A0) = Cx (cos (pi / &2)) where A0:real^2
```

Flyspeck Progress



Neural Autoformalization (Wang et al., 2018)

- generate about 1M Latex - Mizar pairs synthetically (quite advanced)
- train neural seq-to-seq translation models (Luong – NMT)
- evaluate on about 100k examples
- many architectures tested, some work much better than others
- very important latest invention: attention in the seq-to-seq models
- more data crucial for neural training
- Recent addition: unsupervised MT methods (Lample et al 2018) – no need for aligned data, improving a lot!
- Type-checking not yet internal (boosting well-typed data externally)

Neural Autoformalization data

Rendered \LaTeX

If $X \subseteq Y \subseteq Z$, then $X \subseteq Z$.

Mizar

$X \subseteq Y \ \& \ Y \subseteq Z$ implies $X \subseteq Z$;

Tokenized Mizar

$X \subseteq Y \ \& \ Y \subseteq Z$ implies $X \subseteq Z$;

\LaTeX

If $\$X \subseteq Y \subseteq Z\$,$ then $\$X \subseteq Z\$.$

Tokenized \LaTeX

If $\$ X \subseteq Y \subseteq Z \$,$ then $\$ X \subseteq Z \$.$

Neural Autoformalization results

Parameter	Final Test Perplexity	Final Test BLEU	Identical Statements (%)	Identical No-overlap (%)
128 Units	3.06	41.1	40121 (38.12%)	6458 (13.43%)
256 Units	1.59	64.2	63433 (60.27%)	19685 (40.92%)
512 Units	1.6	67.9	66361 (63.05%)	21506 (44.71%)
1024 Units	1.51	61.6	69179 (65.73%)	22978 (47.77%)
2048 Units	2.02	60	59637 (56.66%)	16284 (33.85%)

Neural Fun – Performance after Some Training

Rendered
L^AT_EX

Input L^AT_EX

Correct

Snapshot-
1000

Snapshot-
2000

Snapshot-
3000

Snapshot-
4000

Snapshot-
5000

Snapshot-
6000

Snapshot-
7000

Suppose s_8 is convergent and s_7 is convergent . Then $\lim(s_8+s_7) = \lim s_8 + \lim s_7$

```
Suppose  $\{ s_{8} \}$  is convergent and  $\{ s_{7} \}$ 
is convergent . Then  $\lim ( \{ s_{8} \}
+ \{ s_{7} \} ) \mathrel{=} \lim \{ s_{8} \}
+ \lim \{ s_{7} \}$  .
```

```
seq1 is convergent & seq2 is convergent implies
lim ( seq1 + seq2 ) = ( lim seq1 ) + ( lim seq2 ) ;
```

```
x in dom f implies ( x * y ) * ( f | ( x | ( y | ( y | y )
) ) ) = ( x | ( y | ( y | ( y | y ) ) ) ) ;
```

```
seq is summable implies seq is summable ;
```

```
seq is convergent & lim seq = 0c implies seq = seq ;
```

```
seq is convergent & lim seq = lim seq implies seq1 + seq2
is convergent ;
```

```
seq1 is convergent & lim seq2 = lim seq2 implies lim_inf
seq1 = lim_inf seq2 ;
```

```
seq is convergent & lim seq = lim seq implies seq1 + seq2
is convergent ;
```

```
seq is convergent & seq9 is convergent implies
lim ( seq + seq9 ) = ( lim seq ) + ( lim seq9 ) ;
```

Unsupervised NMT Fun on Short Formulas

```
len <* a *> = 1 ;
assume i < len q ;
len <* q *> = 1 ;
s = apply ( v2 , v1 ast t ) ;
s . ( i + 1 ) = tt . ( i + 1 )
1 + j <= len v2 ;
1 + j + 0 <= len v2 + 1 ;
let i be Nat ;
assume v is_applicable_to t ;
let t be type of T ;
a ast t in downarrow t ;
t9 in types a ;
a ast t <= t ;
A is_applicable_to t ;
Carrier ( f ) c= B
u in B or u in { v } ;
F . w in w & F . w in I ;
GG . y in rng HH ;
a * L = Z_ZeroLC ( V ) ;
not u in { v } ;
u <> v ;
v - w = v1 - w1 ;
v + w = v1 + w1 ;
x in A & y in A ;

len <* a *> = 1 ;
i < len q ;
len <* q *> = 1 ;
s = apply ( v2 , v1 ) . t ;
s . ( i + 1 ) = tau1 . ( i + 1 )
1 + j <= len v2 ;
1 + j + 0 <= len v2 + 1 ;
i is_at_least_length_of p ;
not v is applicable ;
t is_orientedpath_of v1 , v2 , T ;
a *' in downarrow t ;
t '2 in types a ;
a *' <= t ;
A is applicable ;
support ppf n c= B
u in B or u in { v } ;
F . w in F & F . w in I ;
G0 . y in rng ( H1 ./ . y ) ;
a * L = ZeroLC ( V ) ;
u >> v ;
u <> v ;
vw = v1 - w1 ;
v + w = v1 + w1 ;
assume [ x , y ] in A ;
```

Acknowledgments

- Prague Automated Reasoning Group <http://arg.ciirc.cvut.cz/>:
 - Jan Jakubuv, Chad Brown, Martin Suda, Karel Chvalovsky, Bob Veroff, Zar Goertzel, Bartosz Piotrowski, Lasse Blaauwbroek, Martin Smolik, Jiri Vyskocil, Petr Pudlak, David Stanovsky, Krystof Hoder, ...
- HOL(y)Hammer group in Innsbruck:
 - Cezary Kaliszyk, Thibault Gauthier, Michael Faerber, Yutaka Nagashima, Shawn Wang
- ATP and ITP people:
 - Stephan Schulz, Geoff Sutcliffe, Andrej Voronkov, Kostya Korovin, Larry Paulson, Jasmin Blanchette, John Harrison, Tom Hales, Tobias Nipkow, Andrzej Trybulec, Piotr Rudnicki, Adam Pease, ...
- Learning2Reason people at Radboud University Nijmegen:
 - Herman Geuvers, Tom Heskes, Daniel Kuehlwein, Evgeni Tsivtsivadze,
- Google Research: Christian Szegedy, Geoffrey Irving, Alex Alemi, Francois Chollet, Sarah Loos
- ... and many more ...
- Funding: Marie-Curie, NWO, ERC

Some References

- ARG ML&R course: <http://arg.ciirc.cvut.cz/teaching/mlr19/index.html>
- C. Kaliszyk: <http://cl-informatik.uibk.ac.at/teaching/ss18/mltp/content.php>
- C. Kaliszyk, J. Urban, H. Michalewski, M. Olsak: Reinforcement Learning of Theorem Proving. CoRR abs/1805.07563 (2018)
- Z. Goertzel, J. Jakubuv, S. Schulz, J. Urban: ProofWatch: Watchlist Guidance for Large Theories in E. CoRR abs/1802.04007 (2018)
- T. Gauthier, C. Kaliszyk, J. Urban, R. Kumar, M. Norrish: Learning to Prove with Tactics. CoRR abs/1804.00596 (2018).
- J. Jakubuv, J. Urban: ENIGMA: Efficient Learning-Based Inference Guiding Machine. CICM 2017: 292-302
- S. M. Loos, G. Irving, C. Szegedy, C. Kaliszyk: Deep Network Guided Proof Search. LPAR 2017: 85-105
- L. Czajka, C. Kaliszyk: Hammer for Coq: Automation for Dependent Type Theory. J. Autom. Reasoning 61(1-4): 423-453 (2018)
- J. C. Blanchette, C. Kaliszyk, L. C. Paulson, J. Urban: Hammering towards QED. J. Formalized Reasoning 9(1): 101-148 (2016)
- G. Irving, C. Szegedy, A. Alemi, N. Eén, F. Chollet, J. Urban: DeepMath - Deep Sequence Models for Premise Selection. NIPS 2016: 2235-2243
- C. Kaliszyk, J. Urban, J. Vyskocil: Efficient Semantic Features for Automated Reasoning over Large Theories. IJCAI 2015: 3084-3090
- J. Urban, G. Sutcliffe, P. Pudlák, J. Vyskocil: MaLAREa SG1- Machine Learner for Automated Reasoning with Semantic Guidance. IJCAR 2008: 441-456
- C. Kaliszyk, J. Urban, J. Vyskocil: Automating Formalization by Statistical and Semantic Parsing of Mathematics. ITP 2017: 12-27
- Q. Wang, C. Kaliszyk, J. Urban: First Experiments with Neural Translation of Informal to Formal Mathematics. CoRR abs/1805.06502 (2018)
- J. Urban, J. Vyskocil: Theorem Proving in Large Formal Mathematics as an Emerging AI Field. LNCS 7788, 240-257, 2013.

Thanks and Advertisement

- Thanks for your attention!
- **AITP – Artificial Intelligence and Theorem Proving**
- March 22–27, 2020, Aussois, France, aitp-conference.org
- ATP/ITP/Math vs AI/Machine-Learning people, Computational linguists
- Discussion-oriented and experimental - submit a talk abstract!
- Grown to 80 people in 2019