

Caching and Batching for Neural Clause Selection in Saturation-Based Theorem Proving

Alexander Pluska ✉

TU Wien, Austria

Florian Zuleger ✉

TU Wien, Austria

Abstract

A neural clause selector sits in the hot path of a saturation-based theorem prover: the model is invoked on every generated clause and the prover blocks on each result. Whether the ML guidance pays off under a fixed wall budget depends on two axes, the model’s selection quality and the per-call cost the integration imposes. This paper holds the first axis fixed — a GCN encoder paired with a transformer scorer — and isolates two engineering decisions on the second: whether clause embeddings are cached across iterations, and whether embed requests are fired one at a time (eager) or buffered to the selection boundary (batched). On 16 TPTP problems run to a 256-iteration cap, caching changes per-problem wall time by at least a factor of $\approx 28\times$, and without it the prover routinely exhausts the 60-second wall budget; the batching decision is smaller and workload-dependent (eager wins by 1–40% on this set, where $|\Delta| \leq 20$ new clauses are generated per iteration; on heavier workloads the ordering reverses). The two decisions are independent and orthogonal to the model, so the same trade-off recurs in any encoder/scorer-shaped clause selector.

2012 ACM Subject Classification Theory of computation \rightarrow Automated reasoning; Computing methodologies \rightarrow Neural networks

Keywords and phrases automated theorem proving, saturation, clause selection, neural networks, asynchronous inference

Digital Object Identifier 10.4230/LIPIcs...

Supplementary Material *Software*: <https://github.com/lammdachs/proofatlas>

Funding The authors acknowledge the support of the project VASSAL: “Verification and Analysis for Safety and Security of Applications in Life” funded by the European Union under Horizon Europe WIDERA Coordination and Support Action/Grant Agreement No. 101160022.

1 Introduction

Machine learning guidance for saturation-based theorem provers has progressed from hand-crafted clause features [22, 8] to graph neural network encoders over clause graphs [15, 5]. Across these systems the learned component sits in the prover’s hot path: on each newly generated clause the model is invoked, and the prover cannot make its next selection decision until the result is in. For calculi like superposition, which can generate tens of thousands of clauses per proof attempt, the model is invoked tens of thousands of times before a proof is found.

This puts the ML system between two competing forces. The model is chosen for its *selection quality*: a smarter scorer picks clauses that lead to proofs faster. But every invocation costs wall-clock time, and the prover runs against a fixed wall budget; the more time the model takes per call, the fewer iterations the prover gets to take in the same budget. The metric that matters, proof rate within a fixed time, depends on both axes.

The model’s design fixes one axis. The other, the per-call cost, depends on how the model is *invoked* from the prover loop. This paper isolates two engineering decisions that



© Alexander Pluska and Florian Zuleger;

licensed under Creative Commons License CC-BY 4.0

Leibniz International Proceedings in Informatics

LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

XX:2 Caching and Batching for Neural Clause Selection

44 we have found to sit on that second axis. Both are about scheduling around the model, not
45 about the model itself; both are independent of which encoder or scorer architecture is used.

46 Caching of clause embeddings.

47 A clause’s embedding depends only on the clause, not on the rest of the prover state (Sec-
48 tion 2). When a clause enters the unprocessed set, its embedding can be computed once and
49 cached for all subsequent selection steps in which the clause participates. The alternative
50 is to re-embed every clause in both the unprocessed and processed sets at every selection
51 step; this multiplies the backend workload per iteration by a large factor and is, in our
52 measurements, the larger of the two levers.

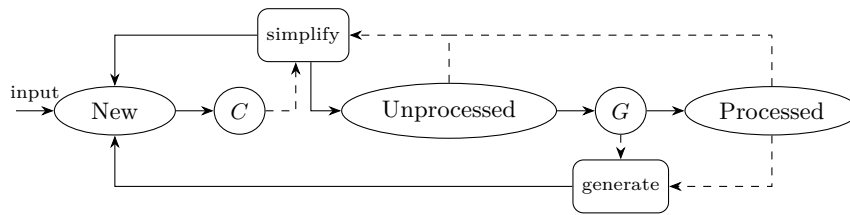
53 Batching at the selection boundary.

54 Within one iteration, the prover’s forward and backward simplification loops generate a set
55 Δ of new clauses; each surviving clause needs an embedding before the next selection step.
56 There is a real choice about *when* the model is invoked to compute those embeddings: either
57 as each new clause arrives, paying any fixed per-call overhead $|\Delta|$ times, or once per iteration
58 over all of Δ together, paying it once. The first pattern lets the model run in parallel with
59 the prover’s continued simplification; the second amortizes a single invocation over more
60 work. Which of these dominates depends on how many clauses an iteration generates and
61 on the relative cost of an invocation.

62 Contributions.

- 63 ■ We isolate the caching decision and the batching decision and measure their independent
64 effects on a single configuration (a GCN encoder paired with a transformer scorer, referred
65 to throughout as *encode_score*) on 16 TPTP problems that the prover runs to its 256-
66 iteration limit under every cell. We run three repetitions per cell.
- 67 ■ We show that the caching decision changes per-problem wall time by at least a factor of
68 $\approx 28\times$ (a lower bound: roughly half of the no-cache runs hit the 60-second wall budget
69 before completing the iteration budget). Without caching the inference cost alone is
70 enough to push these problems past the wall.
- 71 ■ We show that the batching decision is workload-dependent: on the studied iter-cap prob-
72 lems ($|\Delta| \leq 20$ new clauses per iteration) eager submission wins by 1–40% over batched
73 submission, because the embed work overlaps with the prover’s continued simplification;
74 on heavier workloads (e.g. LCL084-10 with $|\Delta| \approx 80$) batched wins instead.
- 75 ■ We place these engineering decisions in their proof-rate context. On a broader 11,091-
76 problem TPTP slice, the cached *encode_score* configuration proves 27.8% of problems
77 versus 26.9% for the non-learned age-weight heuristic. The ML selection quality buys
78 roughly one percentage point of proof rate at roughly $10\times$ the per-problem wall cost of
79 the heuristic, and that comparison only exists because the caching decision keeps the
80 per-problem wall cost from blowing past the time budget.
- 81 ■ We describe the implementation of both decisions in our prover PROOFATLAS¹ and
82 discuss the implications for other ML-guided saturation systems.

¹ Source at <https://github.com/lammdachs/proofatlas>.



■ **Figure 1** The given-clause loop. Dashed arrows show inputs to inference steps; solid arrows show clause flow.

2 Background

2.1 Saturation-Based Theorem Proving

We work in first-order logic with equality. Given a set of clauses S , the goal is to derive the empty clause \square , proving unsatisfiability. To show that a conjecture follows from axioms, we negate the conjecture and prove the combined clause set unsatisfiable.

► **Example 1.** As a running example, we use TPTP problem GRP001-1 [21]: if the square of every element in a group is the identity, the group is commutative. The 11 input clauses include unit clauses such as $product(identity, X, X)$ (left identity), $product(X, Y, multiply(X, Y))$ (defining $multiply$ as a Skolem function for $product$), and $product(X, X, identity)$ (hypothesis), as well as multi-literal clauses such as

$$\neg product(X, Y, Z) \vee \neg product(X, Y, W) \vee Z \approx W$$

encoding well-definedness of $product$.

The given-clause algorithm.

Saturation-based provers iteratively enumerate consequences of the initial clause set, applying *generating inferences* to derive new clauses and *simplifying inferences* to delete or rewrite existing ones. The search terminates with *proof found* when the empty clause is derived, or *saturated* when no further new clauses can be generated. Implementations vary in how they organise the workspace; the version we describe here, and use in PROOFATLAS, maintains three clause sets — *new*, *unprocessed*, and *processed* [24] — and proceeds (Figure 1):

1. While there are clauses in new:
 - a. Pick a clause C ; if $C = \square$, return *proof found*.
 - b. Apply forward simplification; if C is deleted, repeat.
 - c. Apply backward simplification using C ; move C to unprocessed.
2. If unprocessed is empty, return *saturated*. Otherwise, select a clause G (the *given clause*) and move it to processed.
3. Apply generating inferences between G and processed (results go to new).

Clause selection.

The choice of G in step 2—*clause selection*—is the main branching decision of the algorithm. The standard non-learned heuristic is *age-weight interleaving* [17]: selection alternates between the oldest clause (ensuring fairness) and the lightest clause (preferring simpler terms) according to a configurable ratio. Replacing this heuristic with a learned selector is the setting in which the present paper’s contribution operates.

XX:4 Caching and Batching for Neural Clause Selection

115 2.2 The Superposition Calculus

116 The superposition calculus [2] extends resolution with equality handling. Inferences are
117 parametrized by a *term ordering* \succ and a *literal selection strategy*: a function assigning to
118 each clause a subset of its literals.

119 **Generating inferences.**

Underlined literals must be selected. $L[s]$ denotes literal L containing subterm s .

$$\frac{C \vee \underline{L} \quad D \vee \neg \underline{L'}}{(C \vee D)\sigma} \text{ RES} \qquad \frac{C \vee \underline{L} \vee L'}{(C \vee L)\sigma} \text{ FAC}$$

where $\sigma = \text{mgu}(L, L')$ is the most general unifier of L and L' .

$$\frac{l \approx r \vee C \quad L[l'] \vee D}{(L[r] \vee D \vee C)\sigma} \text{ SUP}$$

where $\sigma = \text{mgu}(l, l')$, l' is not a variable, and $l\sigma \not\approx r\sigma$. If $L[l']$ is an equality literal $s[l'] \approx t$ or $s[l'] \not\approx t$, additionally $s[l']\sigma \not\approx t\sigma$.

$$\frac{C \vee s \not\approx t}{C\sigma} \text{ EQRES} \qquad \frac{s \approx s' \vee t \approx t' \vee C}{(s' \not\approx t' \vee s \approx s' \vee C)\sigma} \text{ EQFAC}$$

120 where $\sigma = \text{mgu}(s, t)$, and in EQFAC additionally $s\sigma \not\approx s'\sigma$, $s\sigma \not\approx t'\sigma$, $s'\sigma \not\approx t'\sigma$.

121 **Simplifying inferences.**

122 Four simplifying inferences reduce the search space:

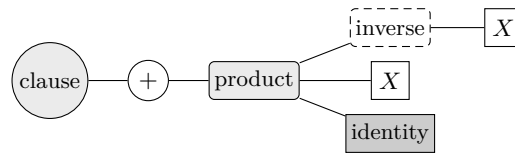
- 123 ■ **Tautology deletion:** Delete a clause that contains complementary literals or a reflexive
124 equality $t \approx t$.
- 125 ■ **Demodulation:** Given a unit equality $l \approx r$ with $l \succ r$ and a substitution σ , replace all
126 occurrences of $l\sigma$ in a clause by $r\sigma$.
- 127 ■ **Condensation:** Replace clause C by $C\sigma$ when there exists a substitution σ such that
128 $C\sigma$ has strictly fewer literals than C and $C\sigma$ subsumes C .
- 129 ■ **Subsumption:** Delete clause D if there exists a clause C and substitution σ such that
130 every literal of $C\sigma$ is a literal of D .

131 **Term ordering.**

132 We use the Knuth–Bendix ordering (KBO) [12] with alphabetic symbol precedence and
133 constant 1 weight function.

134 **Literal selection.**

135 We use the Sel21 strategy of Hoder et al. [7]: select the unique maximal literal if it exists,
136 else a maximal negative literal, else all maximal literals.



■ **Figure 2** Clause graph for $product(inverse(X), X, identity)$. Each symbol occurrence becomes a separate node. Node types: clause root (filled circle), literal (circle, +/−), predicate (filled rounded), function (dashed rounded), variable (plain), constant (dark filled). Edges are bidirectional.

137 2.3 Graph Neural Networks

138 Graph neural networks (GNNs) learn representations of graph-structured data through iter-
 139 ated message passing: each layer updates node representations by aggregating information
 140 from neighbors, and a graph-level representation is obtained by pooling over the final node
 141 embeddings. A graph convolutional network (GCN) [11] updates all node representations
 142 simultaneously by averaging neighbor features with symmetric degree normalization, pre-
 143 venting high-degree nodes from dominating. We apply GraphNorm [3] after each layer,
 144 which normalizes within each graph independently using a learnable mean-subtraction pa-
 145 rameter. A graph-level embedding is obtained by pooling over node representations. GNNs
 146 are widely used as clause encoders in saturation-based theorem proving [14, 8, 6, 5].

147 2.4 Transformers

148 Transformers [23] process sequences through layers of multi-head attention and feed-forward
 149 networks. In scaled dot-product attention, one sequence produces queries while another
 150 produces keys and values; each query computes a dot-product similarity with all keys and
 151 uses the resulting weights to aggregate values. When queries and keys come from the same
 152 sequence, this is *self-attention*; when from different sequences, *cross-attention*. A transformer
 153 layer combines multi-head attention with a position-wise feed-forward network, connected
 154 by residual connections and pre-norm LayerNorm [25].

155 2.5 Learned Clause Selection

156 A learned selector replaces age-weight scoring with a neural model: each unprocessed clause
 157 is encoded into a fixed-dimensional embedding, a scoring head produces a real-valued score
 158 per clause, and the given clause is sampled with temperature-softmax over scores. Encoder
 159 and scorer together form an *encoder-scorer* architecture following TRAIL [5]. We refer to
 160 a specific encoder-scorer pair as a *configuration*. The question of which configuration to
 161 use is largely orthogonal to our contribution; we fix a single configuration — a GCN clause
 162 encoder paired with a transformer scorer (cross-attention from U to P followed by per-token
 163 feed-forward sublayers), which we denote *encode_score* — for all experiments in this paper.
 164 The architectural choice itself, and how the model is trained, are not the subject of this
 165 paper; the questions we study are about how a configuration of this shape should be *invoked*
 166 from inside the prover loop.

167 Clause graphs.

168 The GCN encoder represents each clause as a tree-structured graph (Figure 2) with nodes
 169 for the clause root, literals, predicates, functions, variables, and constants, connected by
 170 bidirectional parent-child edges following the term structure. Each occurrence of a symbol

XX:6 Caching and Batching for Neural Clause Selection

171 or variable becomes a separate node. Each node carries only *structural* features—node
172 type (one-hot), arity, and argument position (sinusoidally encoded). Symbol identity is
173 not used as an input, matching the symbol-independent design of prior GNN-based clause
174 encoders [14, 8]. Importantly, separate clauses constitute separate graphs: there are no
175 edges between clauses, so the encoder treats each clause as an independent input. This
176 makes the per-clause embedding a function of the clause alone, which is what permits the
177 caching strategy of Section 3.

178 2.6 Deployment

179 In PROOFATLAS, neural models run on a single *backend worker* thread with exclusive access
180 to compute resources (CPU threads for small models, one GPU for larger ones). A prover
181 worker submits inference requests to the backend via a channel. The backend’s receive loop
182 drains any concurrently arriving requests (via non-blocking `try_recv`) and dispatches them
183 as a batch to the model.

184 3 Two Implementation Decisions

185 Integrating a learned scorer into the saturation loop forces two implementation decisions. We
186 name them here and measure their effects in Section 5. The two decisions are independent:
187 any combination of choices is a valid implementation.

188 Decision 1: caching of clause embeddings.

189 A clause’s embedding is a function of the clause alone — the encoder consumes the clause
190 graph in isolation, and separate clauses constitute separate graphs (Section 2). Once com-
191 puted at transfer time, the embedding does not need to be re-computed at subsequent
192 selection steps in which the same clause participates. The implementation choice is whether
193 to *cache* the embedding or to recompute it: caching pays the forward-pass cost once per
194 clause and reads it back at every selection step, while recomputing re-embeds every clause
195 in $U \cup P$ at each selection step — $O(|U| + |P|)$ embed calls per iteration rather than $O(|\Delta|)$.
196 Both implementations produce the same scores; they differ only in where the work is done.

197 Decision 2: dispatch timing of embed requests.

198 Within one iteration the prover’s forward and backward simplification loops generate $|\Delta|$
199 new clauses, each needing an embedding before the next selection step. The implementation
200 choice is whether to dispatch each embed request as the corresponding clause is generated
201 (*eager*), or to buffer the $|\Delta|$ requests at the processor and fire them as a single burst at
202 the selection boundary (*batched*). The two patterns trade off in opposite directions: eager
203 submission lets the embed work proceed in parallel with the prover’s continued simplification
204 of subsequent clauses, while batched submission amortizes any fixed per-call overhead (kernel
205 launch, thread wake-up) over the whole iteration’s worth of embeds. Both implementations
206 produce the same scores; they differ only in dispatch schedule.

207 Summary

208 Each of the four combinations (cache vs. no cache, eager vs. batched) is a valid implemen-
209 tation that produces the same proof search modulo floating-point reproducibility. Section 5
210 measures the per-problem wall-time cost of each combination on a 16-problem benchmark.

211 The caching decision turns out to dominate the table by an order of magnitude; the dispatch-
 212 timing decision is smaller and points in different directions depending on the workload.

213 **4 Implementation**

214 We describe how the caching and batching decisions of Section 3 are realized in PROOFAT-
 215 LAS. Both are local changes to the clause-selection processor that sits between the prover’s
 216 `on_transfer/select` callbacks and the inference backend; neither changes the prover, the clause
 217 manager, the saturation calculus, or the model itself.

218 **4.1 Embedding cache**

219 The processor maintains two index maps, $U_{\text{emb}} : \text{idx} \rightarrow \mathbb{R}^d$ and $P_{\text{emb}} : \text{idx} \rightarrow \mathbb{R}^d$, of
 220 clause embeddings for the unprocessed and processed sets respectively. The prover’s clause-
 221 management callbacks update these maps:

- 222 ■ `on_transfer(c)`: record an embed request for c ’s index, populated at the next drain point
 223 (see below) into U_{emb} .
- 224 ■ `on_activate(c)`: move c ’s embedding from U_{emb} to P_{emb} .
- 225 ■ `on_simplify(c)`: remove c ’s embedding from whichever map holds it. If c was simplified
 226 before its embed had fired, drop the pending request entirely.

227 When `select` needs scores, it reads the current U_{emb} and P_{emb} (after firing any pending
 228 embed requests and draining their responses) and submits a single `score_context` call with
 229 the two embedding sets.

230 The cache is a pure performance optimization: a clause’s embedding depends only on
 231 the clause so the cached value is identical to what re-computation would produce. Disabling
 232 the cache (ablated in Section 5) means U_{emb} and P_{emb} are recomputed from scratch at every
 233 `select`, by embedding every clause currently in U and P .

234 **4.2 Batched submit at the selection boundary**

235 The processor maintains a small local buffer $to_submit : \text{Vec}(\langle \text{idx}, \text{PrebuiltInput} \rangle)$ of embed
 236 requests not yet sent to the backend. Each `PrebuiltInput` is the per-clause graph and fea-
 237 ture vector produced by translating the clause into the backend’s input representation; this
 238 translation is done eagerly on the processor thread in `on_transfer`, so that the CPU graph-
 239 building work overlaps with the prover’s continued simplification loop rather than running
 240 on the backend thread. The buffer therefore carries already-tensorized inputs, not the clause
 241 itself. The buffering policy is parameterized by a single integer B , the *embed batch size*:

- 242 ■ In `on_transfer`, the clause is translated into a `PrebuiltInput` and the pair is appended to
 243 to_submit . If $|to_submit| \geq B$, the buffer is flushed (`fire_buffer`): each entry is sent to
 244 the backend via `submit_async` in rapid succession.
- 245 ■ At the next `select` call the buffer is flushed regardless of whether B has been reached,
 246 the backend responses are awaited, and the resulting embeddings are written into U_{emb}
 247 before scoring proceeds.

248 $B = 1$ is the eager policy: each embed request fires its own backend round trip. B
 249 larger than the typical $|\Delta|$ is the batched policy: the buffer fills only at drain time, and
 250 the backend receives one rapid burst of $|\Delta|$ `submit_async` calls per iteration. Because the
 251 backend’s receive loop drains non-blocking after each `recv`, this burst is observed as a single
 252 batch of size $|\Delta|$ and dispatched as one forward pass.

XX:8 Caching and Batching for Neural Clause Selection

253 We set B to a value much larger than any plausible $|\Delta|$ (specifically 10^6 , well above the
254 largest per-iteration $|\Delta|$ we have ever observed) when the batched policy is selected. The
255 buffer cap is conceptually a function of GPU memory available to the embed model; the
256 same memory budget that already accommodates the `score_context` call’s $|U| + |P|$ inputs
257 comfortably accommodates any realistic embed batch.

258 4.3 What does not change

259 The cache and the batching policy are orthogonal: any combination of $\{\text{cache, no cache}\} \times$
260 $\{\text{eager, batched}\}$ is valid. The 2×2 grid is the experimental design of Section 5. Neither
261 lever changes the model, the weights, the prover, or the calculus. The set of clauses passed
262 to `score_context` at each selection step is identical across all four cells.

263 5 Evaluation

264 5.1 Setup

265 All experiments run PROOFATLAS with one prover worker on a Linux host with an NVIDIA
266 RTX A5000 (24 GB). Problems come from TPTP version 9.0.0. Per-problem prover limits
267 are 60 s of wall time, 256 given-clause iterations, and 64 MB of clause storage. These
268 three limits are chosen so that on the problems of interest the iteration cap is the binding
269 constraint: the wall and memory budgets are slack on the cached configurations and the
270 prover terminates at iteration 256 with status `resource_limit`, so per-problem wall time
271 becomes a direct measurement of inference throughput rather than of whether a proof was
272 found. The term ordering is KBO with constant-1 weight (Section 2) and the literal selection
273 strategy is Sel21 [7]. We use the `encode_score` configuration (GCN encoder + transformer
274 scorer, Appendix A); the weights are unchanged across cells.

275 The problem set consists of 16 problems (listed in Appendix B) for which the prover
276 hits the 256-iteration limit under every cell of the ablation. The problems span 8 TPTP
277 domains; mean $|\Delta|/\text{iteration}$ on this set ranges from 3.7 to 19.9 clauses (median 8.2).

278 The four cells of the 2×2 ablation are $\{\text{cache=T, cache=F}\} \times \{\text{batched, eager}\}$; batched
279 uses embed batch size $B = 10^6$ (i.e., flush only at drain), eager uses $B = 1$ (each `on_transfer`
280 submits immediately). For reference we additionally run `cache=T` with `sequential` submit, in
281 which `on_transfer` drains its own response inline before returning, and the non-learned age-
282 weight heuristic, which performs no inference at all. Each cell is run 3 times per problem;
283 reported numbers are per-problem medians over the 3 reps.

284 5.2 Caching

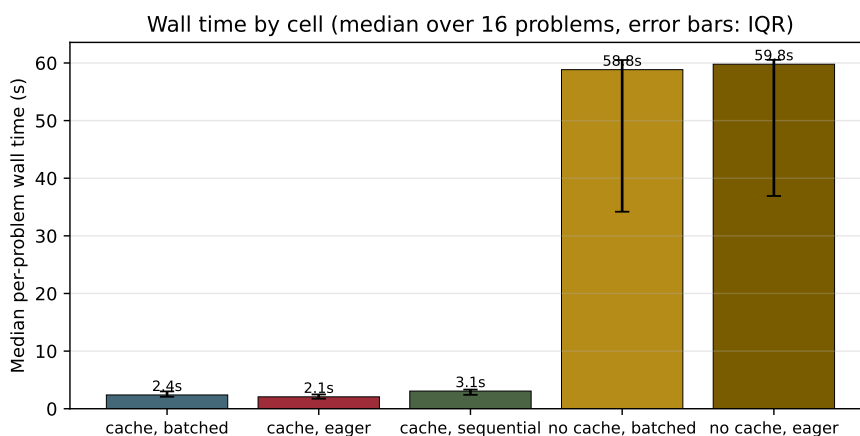
285 Table 1 reports the median per-problem wall time in seconds. The cache vs. no-cache axis
286 dominates the table.

287 The mechanism is direct: without the cache, every clause currently in $U \cup P$ is re-
288 embedded at every `select` call, and $|U|$ grows into the hundreds over the course of the search.
289 The cumulative re-embedding cost dominates the per-problem wall time and routinely ex-
290 ceeds the wall budget — 24 of 48 batched and 25 of 48 eager runs in the `cache=F` row hit the
291 60-second wall before completing the 256-iteration budget. The factor $\approx 28\times$ is consistent
292 across the batched and eager columns, confirming that the caching axis is independent of
293 the dispatch-timing axis.

■ **Table 1** Per-problem wall time (median across 16 problems, 3 reps each). About half of the no-cache runs hit the 60-second wall budget before completing the 256-iteration budget, so the cache=F figures are truncated by the timeout and the 20× ratio is a lower bound.

	batched	eager	sequential
cache=T	2.37 s	2.07 s	3.07 s
cache=F	58.94 s [†]	59.79 s [†]	—
<i>Heuristic baseline (no inference): age-weight</i>			0.22 s

[†] truncated at the 60 s wall budget on roughly half the runs.

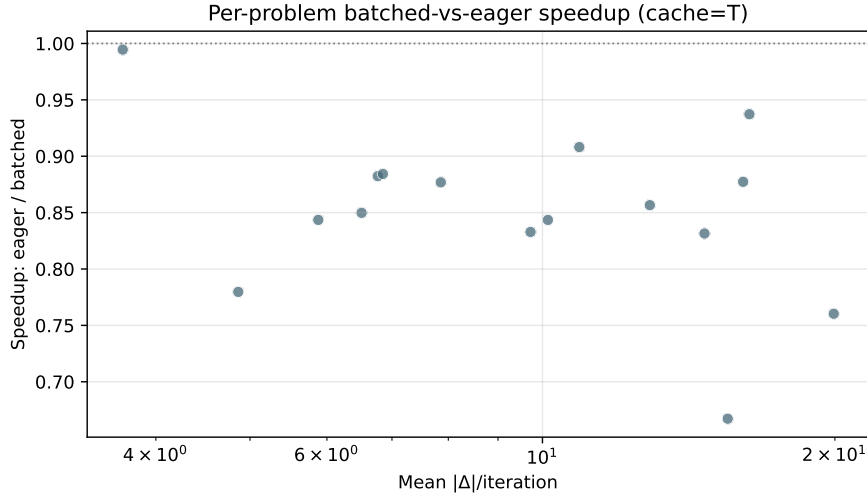


■ **Figure 3** Median per-problem wall time across the 16 problems for each cell of the 2 × 2 ablation plus the sequential-submit reference. Error bars are the inter-quartile range across problems. The two cache=F cells brush against the 60-second wall budget; the ratio between cache rows is a lower bound on the cost of the re-embedding workload.

294 Absolute scale and proof-rate context.

295 The non-learned age-weight heuristic processes the same 16 problems for 256 iterations
 296 each in a median of 0.22 s per problem (single-rep data on the same machine). The cached
 297 encode_score configuration is therefore $\approx 10\times$ slower per problem on the wall-clock axis; the
 298 uncached configuration is roughly two orders of magnitude slower. The encode_score model
 299 does not win on wall time — its per-iteration cost is real — but its *selection decisions* are
 300 better, which is what the ML is there to provide. On an 11,091-problem TPTP slice (all CNF
 301 and FOF problems under 64 kB whose declared status is unsatisfiable or unknown) under
 302 the same 60-second, 256-iteration budget, age-weight proves 2984 problems (26.9 %) and the
 303 cached encode_score configuration proves 3088 problems (27.8 %) — about +1 percentage
 304 point of proof rate that the better selection decisions buy back over the wall-time overhead.
 305 Selection quality and per-iteration cost are the two axes that determine whether an ML-
 306 guided prover beats the heuristic; this paper holds the model (and therefore the selection
 307 quality) fixed and measures what the engineering decisions cost on the second axis. Without
 308 the cache, the per-problem overhead would multiply by another $\approx 20\times$ (per Table 1) and
 309 most TPTP problems would exhaust the 60-second wall before the better selection decisions
 310 had a chance to matter; the modest proof-rate gain would be erased.

XX:10 Caching and Batching for Neural Clause Selection



■ **Figure 4** Per-problem eager/batched wall-time ratio as a function of mean $|\Delta|/\text{iteration}$ across the 16 problems (cache=T). Values below 1.0 mean eager is faster; the dashed line marks parity. The cross-over to batched-faster lies beyond this set’s $|\Delta|$ range.

311 5.3 Batching

312 Within the cache=T row of Table 1, eager (2.07 s) is faster than batched (2.37 s), and
313 sequential is the slowest of the three (3.07 s). The reason: at the low $|\Delta|/\text{iteration}$ values
314 on this problem set (mean 3.7 to 19.9), the prover’s between-call work in the simplify loop
315 is large relative to the cost of one embed forward pass, so firing each embed eagerly lets
316 the backend process it in parallel with that simplify work; by the time select fires, the
317 responses are mostly ready and the drain is fast. The batched policy gives up this overlap,
318 accumulating submits at the processor while the backend is idle. The amortization advantage
319 of batched submission — one kernel launch per iteration instead of $|\Delta|$ — is real but small
320 at these $|\Delta|$ values.

321 Figure 4 shows the per-problem eager/batched wall-time ratio against $|\Delta|/\text{iteration}$; ev-
322 ery problem in this set falls below the parity line, so eager wins everywhere. The ratio
323 trends toward parity as $|\Delta|$ grows, consistent with the prediction that batched submission
324 eventually wins when the per-iteration embed workload is large enough for amortization to
325 exceed the lost overlap. A separate CPU measurement on LCL084-10 (mean $|\Delta| \approx 80$) saw
326 batched faster than eager by $\approx 16\%$, suggesting the cross-over on our setup is somewhere
327 in $|\Delta|/\text{iter} \in [20, 80]$. The right dispatch policy for a given workload therefore depends on
328 which side of that range its typical $|\Delta|$ sits.

329 5.4 Search nondeterminism

330 Even though every cell hits the 256-iteration limit in every run, the four cells process embed
331 requests in batches of different sizes, and deep-learning kernels are not bitwise-stable across
332 batch sizes. Score vectors therefore differ between cells at the level of low bits, which can
333 flip individual selection decisions and lead to different searches. The 256-iteration constraint
334 guarantees that the wall-time numbers across cells reflect throughput on equally long runs,
335 but the specific clauses chosen at each step differ between cells. The wall-time comparisons
336 in this section should be read accordingly.

6 Related Work

6.1 ML-guided saturation provers

Learned clause selection has been explored in E [16] via ENIGMA [9, 10, 8], in Vampire [13, 19, 20], and in TRAIL [5, 1], a saturation prover whose clause-selection policy is itself a neural model with an attention-based action head. These systems typically integrate the learned model as an in-process component invoked from the prover’s main loop. The systems differ in their model architectures (gradient boosting in early ENIGMA variants, graph neural networks over clause graphs in later ENIGMA work, attention over clause embeddings in TRAIL) and in their training data generation strategies, but they share the same overall integration pattern: each new clause must be evaluated before the next selection decision can be made.

Among these systems, the caching and batching decisions we measure are typically not discussed explicitly. ENIGMA’s clause-feature extraction is per-clause, so reusing the extracted feature vector across selection steps is in effect a memoization of the model’s input; the ENIGMA papers do not, however, frame the presence or absence of such reuse as an explicit design decision. Whether the neural-encoder-based saturation selectors [5, 8, 6] cache embeddings or recompute them is, to our knowledge, also not reported, and likewise for whether they batch embed requests across the clauses generated in one iteration.

A recurring theme in this literature is the tension between model expressiveness and inference speed. Suda reports model-inference cost as a significant fraction of overall proof time in Vampire [20]; Jakubův et al. in ENIGMA Anonymous compare gradient-boosted decision trees and graph neural networks head-to-head as clause selectors, and a substantial part of that paper is the engineering required to bring the GNN-based selector’s wall-clock cost into a range where its selection quality is competitive with the cheaper GBDT alternative in real-time proof counts and complementary to it on hard problems [8]. Our work is complementary to this line: rather than making the model cheaper, we characterize how the cost of invoking a fixed model depends on integration choices.

6.2 Clause encoders and selection heads

Recent neural clause-selection systems for saturation provers combine graph neural network encoders over clause graphs [5, 8] with various scoring heads, including the attention-based selection policy of TRAIL [5]. Graph encoders of logical formulas have also been studied outside the saturation-selection setting, for instance for tactic and premise prediction over HOL Light formulas in the HOList benchmark [15]; we reuse the general representation idea but our integration setting and target task are distinct. Our contribution does not propose a new encoder or new scorer: we pair a GCN encoder over clause graphs with a Transformer scorer operating on clause-level embeddings, with cross-attention from the unprocessed set to the processed set (Appendix A). The caching and batching decisions we measure apply identically to other encoder/scorer choices of similar shape.

6.3 Inference serving and batch scheduling

In the broader ML serving literature, request batching at a shared backend is standard practice [4, 18]; systems like Clipper and Nexus expose APIs for the very pattern we describe at the process level. Our contribution is to show that the same pattern matters in the tight-loop in-process setting of a saturation prover — where the “client” (prover) and the “server”

XX:12 Caching and Batching for Neural Clause Selection

380 (inference backend) share an address space and the client issues many requests per iteration
381 — and that getting it wrong has measurable proof-rate consequences.

382 6.4 Parallel saturation

383 Prior work on multi-strategy and parallel saturation (E’s strategy scheduler [17], Vampire’s
384 portfolio modes [13]) explores the *search* axis — running multiple configurations or prover
385 instances against the same problem. Our work is orthogonal: we accelerate a single prover
386 worker’s interaction with the inference backend. How the caching and batching decisions
387 interact with cross-worker scheduling is left for future work.

388 7 Discussion

389 7.1 What this paper is, and is not, about

390 Neither the embedding cache nor the batched-submit policy is novel in any deep sense.
391 Caching of a quantity that depends only on a stable input is how any sensible system handles
392 memoization. Batching of inference requests at a known boundary is how any sensible system
393 handles GPU-side overhead. Yet ML-guided saturation provers in the literature do not
394 typically report which of these decisions they made or measure their impact, so a first-time
395 implementer has no guidance and a re-implementer has no baseline. Our contribution is to
396 isolate the caching and batching decisions, characterize their effects on a single configuration,
397 and demonstrate that each one independently changes per-problem wall time by a factor
398 that matters for the proof-rate metric the field actually cares about.

399 We have intentionally restricted the configuration under study to a single encoder–scorer
400 pair (*encode_score* in the body and Appendix A). The point of the paper is not which
401 encoder or scorer to use; the point is that for any encoder–scorer pair shaped like ours, the
402 cache and the batching policy dominate the throughput question.

403 7.2 Limitations

404 One configuration.

405 We measure one encoder–scorer pair (GCN encoder + transformer scorer). A clause-intrinsic
406 scorer that produces a scalar per clause without consulting a *score_context* over U and P
407 would have a different shape and a different cost profile; in particular, the cache would mem-
408 orize scalars rather than embeddings. Our experimental conclusions are likely qualitatively
409 similar in that case but the per-problem savings should be smaller, because the lookup is
410 cheaper.

411 Single prover worker.

412 All experiments use one prover thread. The natural follow-up is multi-worker scaling, where
413 multiple prover threads share one inference backend. We have intentionally restricted our-
414 selves to a single worker so that the throughput effects we report are attributable to the
415 cache and the batching policy alone, not to scheduling at a shared backend. A proper
416 multi-worker study — with controls for backend queue depth, batch dispatch policy, and
417 floating-point determinism — is left for future work.

418 **Search nondeterminism across cells.**

419 Even when problems hit the iteration limit under every cell of our ablation, the batched
420 and eager policies process embed requests in batches of different sizes, and deep-learning
421 kernels are not bitwise-stable across batch sizes. Score vectors therefore differ at the level
422 of low bits between cells, which can flip individual selection decisions and lead to different
423 searches. We chose problems where this divergence does not flip the proof outcome (Sec-
424 tion 5), but the wall-time numbers across cells should be understood as comparing different
425 search trajectories, not bit-identical ones.

426 **7.3 Future work**

427 **Selection quality versus inference time.**

428 The deepest open question this paper opens is the trade-off the paper itself sidesteps by
429 holding the model fixed: *how much selection quality is worth how much inference time?*
430 A larger, more accurate scorer makes better choices per call but takes longer per call, so
431 under a fixed wall budget it makes fewer choices. The question of which encoder-scorer
432 pair maximises proof rate is therefore not "which is most accurate"; it is "which sits at the
433 right point on a Pareto curve between per-call quality and per-call cost." Characterising that
434 curve, on a fixed proving setup and across model sizes / architectures, is the natural next
435 study and would put the engineering decisions in this paper into context: caching shifts the
436 cost axis (the same model gets much cheaper per problem), batching shifts it more modestly,
437 and together they enlarge the region of model-space that is *affordable* under a given wall
438 budget.

439 **Multi-worker scaling.**

440 With multiple prover workers feeding one inference backend, the dispatch policy at the
441 backend interacts with the within-worker batching policy in non-obvious ways. A controlled
442 study would treat the within- and across-worker batching decisions as orthogonal axes; we
443 expect cross-worker concurrency to push the crossover between eager and batched leftward
444 (toward smaller $|\Delta|$) because the backend already sees more concurrent requests at any given
445 moment.

446 **7.4 Conclusion**

447 A naive ML-guided saturation prover can pay unnecessary cost in two distinct ways: by
448 re-embedding clauses whose embeddings have already been computed, and by dispatching
449 embed requests one at a time when they could be batched together. Each of these is a
450 separate engineering decision; each is independently measurable. In a system of our shape,
451 the cache decision is decisive: caching saves at least a factor of $\approx 28\times$ on the wall-clock cost
452 of running the prover's iteration budget, and without it the prover routinely exhausts the
453 wall budget before completing that budget. The batching decision is smaller and workload-
454 dependent: on problems with few clauses generated per iteration, firing embed requests
455 eagerly is faster because the backend's work overlaps with the prover's simplification; on
456 heavier problems, buffering the requests and firing them as a batch wins because the per-
457 call dispatch overhead gets amortized. We hope that naming these decisions explicitly, and
458 quantifying their impact on a controlled set of problems, helps both new implementers of
459 ML-guided provers and authors of empirical comparisons between them.

AI Disclosure

Claude (Anthropic) was used extensively as a programming assistant in the implementation of PROOFATLAS and for editorial assistance in drafting this paper. All technical contributions, experimental design, and scientific conclusions are the authors' own.

References

- 1 Ibrahim Abdelaziz, Maxwell Crouse, Bassem Makni, Vernon Austel, Cristina Cornelio, Shajith Iqbal, Pavan Kapanipathi, Ndivhuwo Makondo, Kartik Srinivas, Michael Witbrock, and Achille Fokoue. Learning to guide a saturation-based theorem prover. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- 2 Leo Bachmair and Harald Ganzinger. Resolution theorem proving. In *Handbook of Automated Reasoning*, pages 19–99. Elsevier, 2001.
- 3 Tianle Cai, Shengjie Luo, Keyulu Xu, Di He, Tie-yan Liu, and Liwei Wang. GraphNorm: A principled approach to accelerating graph neural network training. In *ICML*, pages 1204–1215, 2021.
- 4 Daniel Crankshaw, Xin Wang, Guilio Zhou, Michael J Franklin, Joseph E Gonzalez, and Ion Stoica. Clipper: A low-latency online prediction serving system. In *NSDI*, 2017.
- 5 Maxwell Crouse, Ibrahim Abdelaziz, Bassem Makni, Spencer Whitehead, Cristina Cornelio, Pavan Kapanipathi, Kartik Srinivas, Veronika Thost, Michael Witbrock, and Achille Fokoue. A deep reinforcement learning approach to first-order logic theorem proving. In *AAAI*, 2021.
- 6 Zarathustra A Goertzel, Karel Chvalovský, Jan Jakubův, Miroslav Olšák, and Josef Urban. Fast and slow enigmas and parental guidance. In *FroCoS*, pages 173–191, 2021.
- 7 Kryštof Hoder, Giles Regeer, Martin Suda, and Andrei Voronkov. Selecting the selection. In *IJCAR*, pages 313–329, 2016.
- 8 Jan Jakubův, Karel Chvalovský, Miroslav Olšák, Bartosz Piotrowski, Martin Suda, and Josef Urban. ENIGMA anonymous: symbol-independent inference guiding machine (system description). In *IJCAR*, 2020.
- 9 Jan Jakubův and Josef Urban. ENIGMA: efficient learning-based inference guiding machine. In *CICM*, volume 10383 of *LNCS*, pages 292–302, 2017.
- 10 Jan Jakubův and Josef Urban. Hammering mizar by learning clause guidance (short paper). In *ITP*, 2019.
- 11 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.
- 12 Donald E Knuth and Peter B Bendix. Simple word problems in universal algebras. pages 263–297, 1970.
- 13 Laura Kovács and Andrei Voronkov. First-order theorem proving and Vampire. In *CAV*, pages 1–35, 2013.
- 14 Miroslav Olšák, Cezary Kaliszyk, and Josef Urban. Property invariant embedding for automated reasoning. In *ECAI*, pages 1395–1402, 2020.
- 15 Aditya Paliwal, Sarah Loos, Markus Rabe, Kshitij Bansal, and Christian Szegedy. Graph representations for higher-order logic and theorem proving. In *AAAI*, 2020.
- 16 Stephan Schulz. E—a brainiac theorem prover. In *AI Communications*, volume 15, pages 111–126, 2002.
- 17 Stephan Schulz. System description: E 1.8. In *LPAR*, pages 735–743, 2013.
- 18 Haichen Shen, Lequn Chen, Yuchen Jin, Liangyu Zhao, Bingyu Kong, Matthai Philipose, Arvind Krishnamurthy, and Ravi Sundaram. Nexus: a GPU cluster engine for accelerating DNN-based video analysis. In *SOSP*, 2019.
- 19 Martin Suda. Improving ENIGMA-style clause selection while learning from history. In *CADE*, 2021.
- 20 Martin Suda. Vampire with a brain is a good ITP hammer. In *FroCoS*, 2021.

- 509 21 Geoff Sutcliffe. The tptp problem library and associated infrastructure: From cnf to th0, tptp
510 v6.4.0. volume 59, pages 483–502, 2017.
- 511 22 Josef Urban, Geoff Sutcliffe, Petr Pudlák, and Jiří Vyskočil. MaLARea SG1 — Machine
512 Learner for Automated Reasoning with Semantic Guidance. In *IJCAR*, volume 5195 of *LNCS*,
513 2008.
- 514 23 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
515 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008,
516 2017.
- 517 24 Uwe Waldmann, Sophie Touret, Simon Robillard, and Jasmin Blanchette. A comprehensive
518 framework for saturation theorem proving. In *IJCAR*, pages 316–334, 2020.
- 519 25 Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang,
520 Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer archi-
521 tecture. In *ICML*, pages 10524–10533, 2020.

522 **A** Encoder–scorer architecture

523 This appendix documents the `encode_score` configuration used in all experiments in this
524 paper. Hyperparameters not listed below take the framework’s defaults.

525 **GCN clause encoder.**

526 Three graph-convolutional layers [11] with hidden dimension $d = 128$, ReLU activations, and
527 GraphNorm after each layer [3]. Node features are concatenated from node type (one-hot, 6
528 dim), arity (log-scaled, 1 dim), and argument position (sinusoidal encoding, 8 dim); symbol
529 identity is not used as an input. The features are linearly projected to d before the first
530 GCN layer. Clause-level pooling uses $1/\sqrt{N}$ -scaled sum pooling, where N is the number of
531 nodes in the clause graph.

532 **Transformer scorer.**

533 Four cross-attention layers with $h = 4$ heads, hidden dimension d , and pre-norm Layer-
534 Norm [25]. Unprocessed clause embeddings serve as queries; processed clauses provide keys
535 and values. Each layer contains an attention sub-layer and a position-wise feed-forward
536 sub-layer (hidden multiplier 4). The final clause representation is projected to a scalar
537 score.

538 **Selection.**

539 At each given-clause step, scores over the unprocessed set are mapped to a categorical
540 distribution by softmax with temperature $\tau = 1$, and one clause is sampled.

541 **B** Throughput-ablation problem set

542 The 16 TPTP problems used in Section 5 are, by domain: GEO082+1, GEO492+1 (geome-
543 try); GRP317-1, GRP366-1 (groups); LCL337-3, LCL670+1.010 (logic calculi); NUM255-1,
544 NUM543+1 (number theory); SET017-4, SET175-6 (set theory); SEU218+3, SEU425+1
545 (set theory in untyped form); SWC063-1, SWC310-1 (software certification); SYN208-1,
546 SYN683-1 (syntactic). All sixteen hit the 256-iteration cap under every cell of the ablation;
547 the per-problem mean $|\Delta|$ per iteration is reported in the project’s experiment log.