

Panbench: A Comparative Benchmarking Tool for Dependently-Typed Languages

3 **Anonymous author**

4 Anonymous affiliation

5 **Anonymous author**

6 Anonymous affiliation

7 — Abstract —

8 We benchmark four proof assistants (Agda, Idris 2, Lean 4 and Rocq) through a single test suite.
9 We focus our benchmarks on the basic features that all systems based on a similar foundations
10 (dependent type theory) have in common. We do this by creating an “over language” in which to
11 express all the information we need to be able to output *correct and idiomatic syntax* for each of our
12 targets. Our benchmarks further focus on “basic engineering” of these systems: how do they handle
13 long identifiers, long lines, large records, large data declarations, and so on.

14 Our benchmarks reveals both flaws and successes in all systems. We give a thorough analysis of
15 the results.

16 We also detail the design of our extensible system. It is designed so that additional tests and
17 additional system versions can easily be added. A side effect of this work is a better understanding
18 of the common abstract syntactic structures of all four systems.

19 **2012 ACM Subject Classification** Replace ccsdesc macro with valid one

20 **Keywords and phrases** Add keywords

21 **Digital Object Identifier** 10.4230/LIPIcs.CVIT.2016.23

22 **1 Introduction**

23 Production-grade implementations of dependently typed programming languages are complica-
24 ted pieces of software that implement many intricate and potentially expensive algorithms.
25 As such, large amounts of engineering effort has been dedicated to optimizing these com-
26 ponents. Unfortunately, engineering time is a finite resource, and this necessarily means
27 that other parts of these systems get comparatively less attention. This often results in
28 easy-to-miss performance problems: we have heard anecdotes from a proof assistant developer
29 that a naïve $O(n^2)$ fresh name generation algorithm used for pretty-printing resulted in 100x
30 slowdowns in some pathological cases.

31 This suggests that a benchmarking suite that focuses on these simpler components
32 could reveal some (comparatively) easy potential performance gains. Moreover, such a
33 benchmarking suite would also be valuable for developers of new dependently typed languages,
34 as it is much easier to optimize with a performance goal in mind. This is an instance of
35 the classic $m \times n$ language tooling problem: constructing a suite of m benchmarks for n
36 languages directly requires a quadratic amount of work up front, and adding either a new
37 test case or a new language to the suite requires an additional linear amount of effort.

38 Like most $m \times n$ tooling problems, the solution is to introduce a mediating tool. In our
39 case, we ought to write all of the benchmarks in an intermediate language, and then translate
40 that intermediate language to the target languages in question. There are existing languages
41 like Dedukti[2] or Informath[1] that attempt to act as an intermediary between popular proof
42 assistants, but these tools typically focus on translating the *content* of proofs, not exact
43 syntactic structure. To fill this gap, we have created the Panbench system, which consists of:



© Anonymous author(s);

licensed under Creative Commons License CC-BY 4.0

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:9

Leibniz International Proceedings in Informatics



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

23:2 Panbench: A Comparative Benchmarking Tool for Dependently-Typed Languages

- 44 1. An extensible embedded Haskell DSL that encodes the concrete syntax of a typical
45 dependently typed language.
- 46 2. A series of compilers for that DSL to Agda, Idris 2, Lean 4, and Rocq.
- 47 3. A benchmarking harness that can perform sandboxed builds of multiple revisions Agda,
48 Idris 2, Lean 4, Rocq.
- 49 4. An incremental build system that can produce benchmarking reports as static HTML
50 files or PGF plots¹.

51 2 Methodology

52 This is really documenting the 'experiment'. The actual details of the thinking behind
the design is in Section 6.

53 this itemized
list should be
expanded into
actual text

- 54 └─ single language of tests
55 └─ document the setup of tests, high level
56 └─ document the setup of testing infrastructure, high level
57 └─ linear / exponential scaling up of test 'size'

58 3 Results

59 4 Discussion

60 5 Infrastructure

61 One of the major goals of Panbench is to make performance analysis as low-cost as possible for
62 language developers. Meeting this goal requires a large amount of supporting infrastructure:
63 simply generating benchmarks is not very useful if you cannot run them nor analyze their
64 outcomes. After some discussion, we concluded that any successful language benchmarking
65 system should meet the following criteria:

- 66 1. It must provide infrastructure for performing sandboxed builds of compilers from source.
67 Asking potential users to set up four different toolchains presents an extremely large
68 barrier to adoption. Moreover, if we rely on user-provided binaries, then we have no hope
69 of obtaining reproducible results, which in turn makes any insights far less actionable.
- 70 2. It must allow for multiple revisions of the same tool to be installed simultaneously. This
71 enables developers to easily look for performance regressions, and quantify the impact of
72 optimizations.
- 73 3. It must allow for multiple copies of the *same* version tool to be installed with different
74 build configurations. This allows developers to look for performance regressions induced
75 by different compiler versions/optimizations.
- 76 4. It must be able to be run locally on a developer's machine. Cloud-based tools are often
77 cumbersome to use and debug, which in turn lowers adoption.
- 78 5. It must present a declarative interface for creating benchmarking environments and
79 running benchmarks. Sandboxed builds of tools are somewhat moot if we cannot trust
80 that a benchmark was run with the correct configuration.

¹ All plots in this paper were produced directly by Panbench!

81 6. It must present performance results in a self-contained format that is easy to understand
82 and share. Performance statistics that require large amounts of post-processing or
83 dedicated tools to view can not be easily shared with developers, which in turn makes
84 the data less actionable.

85 Of these criteria, the first four present the largest engineering challenge, and are tan-
86 tamount to constructing a build system that is able orchestrate *other* build systems. We
87 approached the problem by constructing a bespoke content-addressed system atop of Shake [8],
88 which we discuss in section 5.1. The final two criteria also presented some unforseen difficulties,
89 which we detail in 5.2.

90 5.1 The Panbench Build System

91 As noted earlier, we strongly believe that any benchmarking system should provide infra-
92 structure for installing multiple versions of reproducibly built software. Initially, we intended
93 to build this infrastructure for Panbench atop of Nix [4]. This is seemingly a perfect fit;
94 after all, Nix was designed to facilitate almost exactly this use-case. However, after further
95 deliberation, we came to the conclusion that Nix did not quite meet our needs for the
96 following reasons:

- 97 1. Nix does not work natively on Windows. Performance problems can be operating system
98 specific, so ruling out an OS that has a large user base that is often overlooked in testing
99 seems unwise.
- 100 2. Nix adds a barrier to adoption. Installing Nix is a somewhat invasive process, especially
101 on MacOS². We believe that it is somewhat unreasonable to ask developers to add users
102 and modify their root directory to run a benchmarking tool, and strongly suspect that
103 this would hamper adoption.

104 With the obvious option exhausted, we opted to create our own Nix-inspired build system
105 based atop Shake [8].

106 5.2 Running Benchmarks and Generating Reports

107 6 The Design of Panbench

108 At its core, Panbench is a tool for producing grammatically well-formed concrete syntax
109 across multiple different languages. Crucially, Panbench does *not* require that the syntax
110 produced is well-typed or even well-scoped: if it did, then it would be impossible to benchmark
111 how systems perform when they encounter errors. This seemingly places Panbench in stark
112 contrast with other software tools for working with the meta-theoretic properties of type
113 systems, which are typically concerned only with well-typed terms.

114 However, core task of Panbench is not that different from the task of a logical framework [5]:
115 Both systems exist to manipulate judgements, inference rules, and derivations: Panbench
116 just works with *grammatical* judgements and production rules rather than typing judgments
117 and inference rules. In this sense Panbench is a *grammatical* framework³ rather than a logical
118 one.

Reed: Cite something

² The situation is even worse on x86-64 Macs, which most Nix installers simply do not support.

³ Not to be confused with *the Grammatical Framework* [9], which aims to be a more natural-language focused meta-framework for implementing and translating between grammars.

23:4 Panbench: A Comparative Benchmarking Tool for Dependently-Typed Languages

119 This similarity let us build Panbench atop well-understood design principles. In particular,
120 a mechanized logical framework typically consists of two layers:

- 121 1. A layer for defining judgements à la relations.
122 2. A logic programming layer for synthesizing derivations.

123 To use a logical framework, one first encodes a language by laying out all of the judgements.
124 Then, one needs to prove an adequacy theorem on the side that shows that their encoding of
125 the judgements actually aligns with the language. However, if one wanted to mechanize this
126 adequacy proof, then a staged third layer that consists of a more traditional proof assistant
127 would be required.

128 If we take this skeleton of a design and transpose it to work with grammatical constructs
129 rather than logical ones, we will also obtain three layers:

- 130 1. A layer for defining grammars à la relations.
131 2. A logic programming layer for synthesizing derivations.
132 3. A staged functional programming layer for proving “adequacy” results.

133 In this case, an adequacy result for a given language \mathcal{L} is a constructive proof that all
134 grammatical derivations written within the framework can be expressed within the concrete
135 syntax of a language \mathcal{L} . However, the computational content of such a proof essentially
136 amounts to a compiler written in the functional programming layer.

137 6.1 Implementing The Grammatical Framework

138 Implementing a bespoke hybrid of a logic and functional programming language is no small
139 feat, and also requires prospective users to learn yet another single-purpose tool. Luckily,
140 there already exists a popular, industrial-grade hybrid logic/functional programming language
141 in wide use: GHC Haskell.

142 At first glance, Haskell does not contain a logic programming language. However, if we
143 enable enough GHC extensions, the typeclass system can be made *just* rich enough to encode
144 a simply-typed logical framework. The key insight is that we can encode each production
145 rule using multi-parameter type classes with a single method. Moreover, we can encode our
146 constructive adequacy proofs for a given set of production rules as instances that translate
147 each of the productions in the abstract grammar to productions in the syntax of an actual
148 language.

149 As a concrete example, consider the grammar of the following simple imperative language.

```
150 <expr> := x
151   | n
152   | <expr> '+' <expr>
153   | <expr> '*' <expr>
154
155 <stmt> := <var> '=' <expr>
156   | 'while' <expr> 'do' <stmt>
157   | <stmt> ';' <stmt>
```

159 We can then encode this grammar with the following set of multi-parameter typeclasses:

160  **Listing 1** An example tagless encoding.

```
161 class Var expr where
162   var :: String → expr
```

```

163 class Lit expr where
164   lit :: Int → expr
165
166 class Add expr where
167   add :: expr → expr → expr
168
169 class Mul expr where
170   mul :: expr → expr → expr
171
172 class Assign expr stmt where
173   assign :: String → expr → stmt
174
175 class While expr stmt where
176   while :: expr → stmt → stmt
177
178 class AndThen stmt where
179   andThen :: stmt → stmt → stmt
180
181

```

182 This style of language encoding is typically known as the untyped variant of *finally tagless*[3], and is well-known technique. However, our encoding is a slight refinement of
 183 the usual tagless style. In particular, we restrict ourselves to a single class per production
 184 rule, whereas other tagless encodings often use a class per syntactic category. This more
 185 fine-grained approach allows us to encode grammatical constructs that are only supported
 186 by a subset of our target grammars; see section 6.2 for examples.
 187

188 Unfortunately, the encoding above has some serious ergonomic issues. In particular,
 189 expressions like `assign "x"` (`lit 4`) will result in an unsolved metavariable for `expr`, as there
 190 may be multiple choices of `expr` to use when solving the `Assign ?expr stmt` constraint. Luckily,
 191 we can resolve ambiguities of this form through judicious use of functional dependencies[6],
 192 as demonstrated below.

■ Listing 2 A tagless encoding with functional dependencies.

```

193 class Assign expr stmt | stmt → expr where
194   assign :: String → expr → stmt
195
196 class While expr stmt | stmt → expr where
197   while :: expr → stmt → stmt
198
199

```

200 6.2 Designing The Language

201 Now that we've fleshed out how we are going to encode our grammatical framework into
 202 our host language, it's time to design our idealized abstract grammar. All of our target
 203 languages roughly agree on a subset of the grammar of non-binding terms: the main sources
 204 of divergence are binding forms and top-level definitions⁴. This is ultimately unsurprisingly:
 205 dependent type theories are fundamentally theories of binding and substitution, so we would
 206 expect some variation in how our target languages present the core of their underlying
 207 theories.

⁴ As we shall see in section 6.2.2, the syntactical divergence of top-level definitions is essentially about binders as well.

23:6 Panbench: A Comparative Benchmarking Tool for Dependently-Typed Languages

208 This presents an interesting language design problem. Our idealized grammar will need to
209 find some syntactic overlap between all of our four target languages. Additionally, we would
210 also like for our solution to be (reasonably) extensible. Finding the core set of grammatical
211 primitives to accomplish this task is surprisingly tricky, and requires a close analysis of fine
212 structure of binding.

213 6.2.1 Binding Forms

214 As users of dependently typed languages are well aware, a binding form carries much more
215 information than just a variable name and a type. Moreover, this extra information can have
216 a large impact on typechecking performance, as is the case with implicit/visible arguments.
217 To make matters worse, languages often offer multiple syntactic options for writing the same
218 binding form, as is evidenced by the prevalence of multi-binders like $(x\ y\ z : A) \rightarrow B$. Though
219 such binding forms are often equivalent to their single-binder counterparts as *abstract syntax*,
220 they may have different performance characteristics, so we cannot simply lower them to a
221 uniform single-binding representation. To account for these variations, we have designed a
222 sub-language dedicated solely to binding forms. This language classifies the various binding
223 features along three separate axes: binding arity, binding annotations, and binding modifiers.

224 Binding arities and annotations are relatively self-explanatory, and classify the number of
225 names bound, along with the type of annotation allowed. Our target languages all have their
226 binding arities falling into one of three classes: *n*-ary, unary, or nullary. We can similarly
227 characterise annotations into three categories: required, optional, or forbidden.

228 This language of binders is encoded in the implementation as a single class `Binder` that is
229 parameterised by higher-kinded types `arity :: Type -> Type` and `ann :: Type -> Type`, and
230 provide standardized types for all three binding arities and annotations.

231 Listing 3 The basic binding constructs in Panbench.

```
232 class Binder arity nm ann tm cell | cell -> nm tm where
233   binder :: arity nm -> ann tm -> cell
234
235   -- | No annotation or arity.
236   data None nm = None
237
238   -- | A single annotation or singular arity.
239   newtype Single a = Single { unSingle :: a }
240
241   -- | Multi-binders.
242   type Multi = []
243
244   -- | Infix operator for an annotated binder with a single name.
245   (.::) :: (Binder Single nm Single tm cell) => nm -> tm -> cell
246   nm .:: tp = binder (Single nm) (Single tp)
247
248   -- | Infix operator for an annotated binder.
249   (.::* ) :: (Binder arity nm Single tm cell) => arity nm -> tm -> cell
250   nms .::* tp = binder nms (Single tp)
```

252 Production rules that involve binding forms are encoded as classes that are parametric over
253 a notion of a binding cell, as demonstrated below.

254 Listing 4 The Panbench class for II-types.

```
255 class Pi cell tm | tm -> cell where
```

```
256     pi :: [cell] → tm → tm
```

258 Decoupling the grammar of binding forms from the grammar of binders themselves allows
 259 us to be somewhat polymorphic over the language of binding forms when writing generators.
 260 This in turn means that we can potentially re-use generators when extending Panbench with
 261 new target grammars that may support only a subset of the binding features present in our
 262 four target grammars.

263 Binding modifiers, on the other hand, require a bit more explanation. A binding modifier
 264 captures features like implicit arguments, which do not change the number of names bound
 265 nor their annotations, but rather how those bound names get treated by the rest of the
 266 system. Currently, Panbench only supports visibility-related modifiers, but we have designed
 267 the system so that it is easy to extend with new modifiers; EG: quantities in Idris 2 or
 268 irrelevance annotations in Agda.

269 The language of binding modifiers is implemented as the following set of Haskell type-
 270 classes.

Listing 5 Typeclasses for binding modifiers.

```
271 class Implicit cell where
272   implicit :: cell → cell
273
274 class SemiImplicit cell where
275   semiImplicit :: cell → cell
```

278 This is a case where the granular tagless encoding shines. Both Lean 4 and Rocq have a
 279 form of semi-implicits⁵, whereas Idris 2 and Agda have no such notion. Decomposing the
 280 language of binding modifiers into granular pieces lets us write benchmarks that explicitly
 281 require support for features like semi-implicits. Had we used a monolithic class that encodes
 282 the entire language of modifiers, we would have to resort to runtime errors (or, even worse,
 283 dubious attempts at translation).

284 6.2.2 Top-Level Definitions

285 The question of top-level definitions is much thornier, and there seems to be less agreement
 286 on how they ought to be structured. Luckily, we can re-apply many of the lessons we
 287 learned in our treatment of binders; after all, definitions are “just” top-level binding forms!
 288 This perspective lets us simplify how we view some more baroque top-level bindings. As a
 289 contrived example, consider the following signature for a pair of top-level Agda definitions.

Listing 6 A complicated Agda signature.

```
290 private instance abstract @irr @mixed foo bar : Nat → _
```

293 In our language of binders, this definition consists of a 2-ary annotated binding of the names
 294 `foo`, `bar` that has had a sequence of binding modifiers applied to it.

295 Unfortunately, this insight does not offer a complete solution. Notably, our four target
 296 grammar differ significantly in how their treatment of type signatures. prioritize dependent
 297 pattern matching (EG: Agda, Idris 2) typically opt to have standalone type signatures: this
 298 allow for top-level pattern matches, which in turn makes it much easier to infer motives[7].

⁵ We use the term *semi-implicit* argument to refer to an implicit argument that is not eagerly instantiated. In Rocq these are known as non-maximal implicits.

23:8 Panbench: A Comparative Benchmarking Tool for Dependently-Typed Languages

299 Conversely, languages oriented around tactics (EG: Lean 4, Rocq) typically opt for in-line
300 type signatures and pattern-matching expressions. This appears to be largely independent of
301 Miranda/ML syntax split: Lean 4 borrows large parts of its syntax from Haskell, yet still
302 opts for in-line signatures.

303 This presents us with a design decision: should our idealized grammar use inline or
304 standalone signatures? As long as we can (easily) translate from one style to the other, we
305 have a genuine decision to make. We have opted for the former as standalone signatures
306 offer variations that languages with inline signatures cannot handle. As a concrete example,
307 consider the following Agda declaration:

Listing 7 A definition with mismatched names.

```
308 id : (A : Type) → A → A
309 id B x = x
310
```

312 In particular, note that we have bound first argument to a different name. Translating
313 this to a corresponding Rocq declaration then forces us to choose to use either the name from
314 the signature or the term. Conversely, using in-line signatures does not lead us to having to
315 make an unforced choice when translating to a separate signature, as we can simply duplicate
316 the name in both the signature and term.

317 However, inline signatures are not completely without fault, and cause some edge cases
318 with binding modifiers. As an example, consider the following two variants of the identity
319 function in Agda.

Listing 8 Two variants of the identity function.

```
320 id : {A : Type} → A → A
321 id x = x
322
323 id' : {A : Type} → A → A
324 id' {A} x = x
325
```

327 Both definitions mark the `A` argument as an implicit, but the second definition *also* binds
328 it in the declaration. When we pass to inline type signatures, we lose this extra layer of
329 distinction. To account for this, we were forced to refine the visibility modifier system to
330 distinguish between “bound” and “unbound” modifiers. This extra distinction has not proved
331 to be too onerous in practice, and we still believe that inline signatures are the correct choice
332 for our application.

333 We have encoded this decision in our idealized grammar by introducing a notion of a
334 “left-hand-side” of a definition, which consists of a collection of names to be defined, and a
335 scope to define them under. This means that we view definitions like `???` not as functions
336 `id : (A : Type) → A → A` but rather as *bindings* `A : Type, x : A ⊢ id : A` in non-empty
337 contexts. This shift in perspective has the added benefit of making the interface to other
338 forms of parameterised definitions entirely uniform; for instance, a parameterised record is
339 simply just a record with a non-empty left-hand side.

340 In Panbench, definitions and their corresponding left-hand sides are encoded via the
341 following set of typeclasses.

Listing 9 Definitions and left-hand sides.

```
342 class Definition lhs tm defn | defn → lhs tm where
343   (.=) :: lhs → tm → defn
344
345 class DataDefinition lhs ctor defn | defn → lhs ctor where
```

Reed: link to
the previous
listing

```
347     data_ :: lhs → [ctor] → defn
348
349 class RecordDefinition lhs nm fld defn | defn → lhs nm fld where
350   record_ :: lhs → name → [fld] → defn
351
```

352 7 Conclusion

353 References

- 354 1 February 2026. URL: <https://github.com/GrammaticalFramework/informath>.
- 355 2 Ali Assaf, Guillaume Burel, Raphaël Cauderlier, David Delahaye, Gilles Dowek, Catherine Dubois, Frédéric Gilbert, Pierre Halmagrand, Olivier Hermant, and Ronan Saillard. Dedukti: a logical framework based on the $\lambda\pi$ -calculus modulo theory. (arXiv:2311.07185), November 2023. arXiv:2311.07185 [cs]. URL: <http://arxiv.org/abs/2311.07185>, doi:10.48550/arXiv.2311.07185.
- 360 3 Jacques Carette, Oleg Kiselyov, and Chung-Chieh Shan. Finally tagless, partially evaluated: Tagless staged interpreters for simpler typed languages. *Journal of Functional Programming*, 19(5):509–543, 2009. doi:10.1017/S0956796809007205.
- 363 4 Eelco Dolstra and The Nix contributors. Nix. URL: <https://github.com/NixOS/nix>.
- 364 5 Robert Harper, Furio Honsell, and Gordon Plotkin. A framework for defining logics. *Journal of the ACM*, 40(1):143–184, January 1993. doi:10.1145/138027.138060.
- 366 6 Mark P. Jones. Type classes with functional dependencies. In Gert Smolka, editor, *Programming Languages and Systems*, page 230–244, Berlin, Heidelberg, 2000. Springer. doi:10.1007/3-540-46425-5_15.
- 369 7 Conor McBride and James McKinna. The view from the left. *Journal of Functional Programming*, 14(1):69–111, January 2004. doi:10.1017/S0956796803004829.
- 371 8 Neil Mitchell. Shake before building: replacing make with haskell. *SIGPLAN Not.*, 47(9):55–66, 2012. doi:10.1145/2398856.2364538.
- 373 9 Aarne Ranta. *Grammatical framework: programming with multilingual grammars*. CSLI studies in computational linguistics. CSLI Publications, Center for the Study of Language and Information, Stanford, Calif, 2011.