

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 tated 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. 23:1–23:8



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

- Reed: Cita-
tions here
- 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,
Idris 2, Lean 4, Rocq.
- 48 4. An incremental build system that can produce benchmarking reports as static HTML
49 files or PGF plots¹.
- 50

51 2 Methodology

52 This is really documenting the 'experiment'. The actual details of the thinking behind
the design is in Section 5.

this itemized
list should be
expanded into
actual text

- 53
- 54 5. single language of tests
 - 55 5. document the setup of tests, high level
 - 56 5. document the setup of testing infrastructure, high level
 - 57 5. linear / exponential scaling up of test 'size'

58 3 Results

59 4 Discussion

60 5 The Design of Panbench

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

65 However, core task of Panbench is not that different from the task of a logical framework [4]:
66 Both systems exist to manipulate judgements, inference rules, and derivations: Panbench
67 just works with *grammatical* judgements and production rules rather than typing judgments
68 and inference rules. In this sense Panbench is a *grammatical* framework² rather than a logical
69 one.

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

- 72
- 73 1. A layer for defining judgements à la relations.
 - 74 2. A logic programming layer for synthesizing derivations.

75 To use a logical framework, one first encodes a language by laying out all of the judgements.
76 Then, one needs to prove an adequacy theorem on the side that shows that their encoding of
77 the judgements actually aligns with the language. However, if one wanted to mechanize this

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

² Not to be confused with *the Grammatical Framework* [7], which aims to be a more natural-language focused meta-framework for implementing and translating between grammars.

79 adequacy proof, then a staged third layer that consists of a more traditional proof assistant
 80 would be required.

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

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

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

90 5.1 Implementing The Grammatical Framework

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

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

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

```
103 <expr> := x
104   | n
105   | <expr> '+' <expr>
106   | <expr> '*' <expr>
107
108 <stmt> := <var> '=' <expr>
109   | 'while' <expr> 'do' <stmt>
110   | <stmt> ';' <stmt>
```

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

113 ■ **Listing 1** An example tagless encoding.

```
114 class Var expr where
115   var :: String → expr
116
117 class Lit expr where
118   lit :: Int → expr
119
120 class Add expr where
121   add :: expr → expr → expr
122
123 class Mul expr where
124   mul :: expr → expr → expr
```

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

```
126 class Assign expr stmt where
127   assign :: String → expr → stmt
128
129 class While expr stmt where
130   while :: expr → stmt → stmt
131
132 class AndThen stmt where
133   andThen :: stmt → stmt → stmt
```

135 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
136 the usual tagless style. In particular, we restrict ourselves to a single class per production
137 rule, whereas other tagless encodings often use a class per syntactic category. This more
138 fine-grained approach allows us to encode grammatical constructs that are only supported
139 by a subset of our target grammars; see section 5.2 for examples.

141 Unfortunately, the encoding above has some serious ergonomic issues. In particular,
142 expressions like `assign "x"` (lit 4) will result in an unsolved metavariable for `expr`, as there
143 may be multiple choices of `expr` to use when solving the `Assign ?expr stmt` constraint. Luckily,
144 we can resolve ambiguities of this form through judicious use of functional dependencies[5],
145 as demonstrated below.

146 **Listing 2** A tagless encoding with functional dependencies.
147

```
class Assign expr stmt | stmt → expr where
148   assign :: String → expr → stmt
149
150 class While expr stmt | stmt → expr where
151   while :: expr → stmt → stmt
```

153 5.2 Designing The Language

154 Now that we've fleshed out how we are going to encode our grammatical framework into
155 our host language, it's time to design our idealized abstract grammar. All of our target
156 languages roughly agree on a subset of the grammar of non-binding terms: the main sources
157 of divergence are binding forms and top-level definitions³. This is ultimately unsurprisingly:
158 dependent type theories are fundamentally theories of binding and substitution, so we would
159 expect some variation in how our target languages present the core of their underlying
160 theories.

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

166 5.2.1 Binding Forms

167 As users of dependently typed languages are well aware, a binding form carries much more
168 information than just a variable name and a type. Moreover, this extra information can have

³ As we shall see in section 5.2.2, the syntactical divergence of top-level definitions is essentially about binders as well.

169 a large impact on typechecking performance, as is the case with implicit/visible arguments.
 170 To make matters worse, languages often offer multiple syntactic options for writing the same
 171 binding form, as is evidenced by the prevalence of multi-binders like $(x\ y\ z : A) \rightarrow B$. Though
 172 such binding forms are often equivalent to their single-binder counterparts as *abstract* syntax,
 173 they may have different performance characteristics, so we cannot simply lower them to a
 174 uniform single-binding representation. To account for these variations, we have designed a
 175 sub-language dedicated solely to binding forms. This language classifies the various binding
 176 features along three separate axes: binding arity, binding annotations, and binding modifiers.

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

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

Listing 3 The basic binding constructs in Panbench.

```
184 class Binder arity nm ann tm cell | cell -> nm tm where
185   binder :: arity nm -> ann tm -> cell
186
187   -- | No annotation or arity.
188   data None nm = None
189
190   -- | A single annotation or singular arity.
191   newtype Single a = Single { unSingle :: a }
192
193   -- | Multi-binders.
194   type Multi = []
195
196   -- | Infix operator for an annotated binder with a single name.
197   (.::) :: (Binder Single nm Single tm cell) => nm -> tm -> cell
198   nm .: tp = binder (Single nm) (Single tp)
199
200   -- | Infix operator for an annotated binder.
201   (.:*) :: (Binder arity nm Single tm cell) => arity nm -> tm -> cell
202   nms .:* tp = binder nms (Single tp)
203
```

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

Listing 4 The Panbench class for II-types.

```
207 class Pi cell tm | tm -> cell where
208   pi :: [cell] -> tm -> tm
209
```

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

216 Binding modifiers, on the other hand, require a bit more explanation. A binding modifier
 217 captures features like implicit arguments, which do not change the number of names bound
 218 nor their annotations, but rather how those bound names get treated by the rest of the

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

219 system. Currently, Panbench only supports visibility-related modifiers, but we have designed
220 the system so that it is easy to extend with new modifiers; EG: quantities in Idris 2 or
221 irrelevance annotations in Agda.

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

224 **Listing 5** Typeclasses for binding modifiers.

```
225 class Implicit cell where
226   implicit :: cell → cell
227
228 class SemiImplicit cell where
229   semiImplicit :: cell → cell
```

231 This is a case where the granular tagless encoding shines. Both Lean 4 and Rocq have a
232 form of semi-implicits⁴, whereas Idris 2 and Agda have no such notion. Decomposing the
233 language of binding modifiers into granular pieces lets us write benchmarks that explicitly
234 require support for features like semi-implicits. Had we used a monolithic class that encodes
235 the entire language of modifiers, we would have to resort to runtime errors (or, even worse,
236 dubious attempts at translation).

237 5.2.2 Top-Level Definitions

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

243 **Listing 6** A complicated Agda signature.

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

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

248 Unfortunately, this insight does not offer a complete solution. Notably, our four target
249 grammar differ significantly in how their treatment of type signatures. prioritize dependent
250 pattern matching (EG: Agda, Idris 2) typically opt to have standalone type signatures: this
251 allow for top-level pattern matches, which in turn makes it much easier to infer motives[6].
252 Conversely, languages oriented around tactics (EG: Lean 4, Rocq) typically opt for in-line
253 type signatures and pattern-matching expressions. This appears to be largely independent of
254 Miranda/ML syntax split: Lean 4 borrows large parts of its syntax from Haskell, yet still
255 opts for in-line signatures.

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

⁴ 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.

Listing 7 A definition with mismatched names.

```
261   id : (A : Type) → A → A
262   id B x = x
263
```

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

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

Listing 8 Two variants of the identity function.

```
273   id : {A : Type} → A → A
274   id x = x
275
276
277   id' : {A : Type} → A → A
278   id' {A} x = x
```

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

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

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

Reed: link to the previous listing

Listing 9 Definitions and left-hand sides.

```
295 class Definition lhs tm defn | defn → lhs tm where
296   (.=) :: lhs → tm → defn
297
298 class DataDefinition lhs ctor defn | defn → lhs ctor where
299   data_ :: lhs → [ctor] → defn
300
301 class RecordDefinition lhs nm fld defn | defn → lhs nm fld where
302   record_ :: lhs → name → [fld] → defn
303
```

305 6 Conclusion

306 References

307 1 February 2026. URL: <https://github.com/GrammaticalFramework/informath>.

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

- 308 2 Ali Assaf, Guillaume Burel, Raphaël Cauderlier, David Delahaye, Gilles Dowek, Catherine
309 Dubois, Frédéric Gilbert, Pierre Halmagrand, Olivier Hermant, and Ronan Saillard. Dedukti:
310 a logical framework based on the $\lambda\pi$ -calculus modulo theory. (arXiv:2311.07185), November
311 2023. arXiv:2311.07185 [cs]. URL: <http://arxiv.org/abs/2311.07185>, doi:10.48550/arXiv.
312 2311.07185.
- 313 3 Jacques Carette, Oleg Kiselyov, and Chung-Chieh Shan. Finally tagless, partially evaluated:
314 Tagless staged interpreters for simpler typed languages. *Journal of Functional Programming*,
315 19(5):509–543, 2009. doi:10.1017/S0956796809007205.
- 316 4 Robert Harper, Furio Honsell, and Gordon Plotkin. A framework for defining logics. *Journal
317 of the ACM*, 40(1):143–184, January 1993. doi:10.1145/138027.138060.
- 318 5 Mark P. Jones. Type classes with functional dependencies. In Gert Smolka, editor, *Programming
319 Languages and Systems*, page 230–244, Berlin, Heidelberg, 2000. Springer. doi:10.1007/
320 3-540-46425-5_15.
- 321 6 Conor McBride and James McKinna. The view from the left. *Journal of Functional Programming*,
322 14(1):69–111, January 2004. doi:10.1017/S0956796803004829.
- 323 7 Aarne Ranta. *Grammatical framework: programming with multilingual grammars*. CSLI
324 studies in computational linguistics. CSLI Publications, Center for the Study of Language and
325 Information, Stanford, Calif, 2011.