

# 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** Mathematics of computing → Mathematical software performance

20 **Keywords and phrases** Benchmarking, dependent types, testing

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 feature many intricate and potentially expensive algorithms. As  
25 such, large amounts of engineering effort has been dedicated to optimizing these components.  
26 Unfortunately, engineering time is a finite resource, and this necessarily means that other  
27 parts of these systems get comparatively less attention. This often results in easy-to-miss  
28 performance problems: we have heard anecdotes from a proof assistant developer that a naïve  
29  $O(n^2)$  fresh name generation algorithm used for pretty-printing resulted in 100x slowdowns  
30 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 [3] or Informath [1] that attempt to act as an intermediary between popular  
42 proof 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. 23:1–23:11

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 [2], Idris 2 [4], Lean 4 [6], and Rocq [9].
- 47 3. A benchmarking harness that can perform sandboxed builds of multiple versions of Agda,  
48 Idris 2, Lean 4, and Rocq.
- 49 4. An incremental build system that can produce benchmarking reports as static HTML  
50 files or PGF plots<sup>1</sup>.

### 51 2 Methodology

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

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 Given that our test suite has 32 tests, each of which produces 3 different graphs, we have no  
60 room to display all 96 resulting graphs<sup>2</sup>. We thus choose results that appear to be the most  
61 "interesting".

### 62 4 Discussion

63 The previous section analyzed the results themselves. Here we speculate on why the results  
64 may be as they are. We comment on some particular results first, and then on what we find  
65 for each system.

#### 66 4.0.0.1 General

#### 67 4.0.0.2 Agda

#### 68 4.0.0.3 Idris 2

#### 69 4.0.0.4 Lean 4

#### 70 4.0.0.5 Rocq

### 71 5 Infrastructure

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

---

<sup>1</sup> All plots in this paper were produced directly by Panbench.

<sup>2</sup> Nor can we have appendices!

- 77 1. It must provide infrastructure for performing sandboxed builds of compilers from source.  
78 Asking potential users to set up four different toolchains presents an extremely large  
79 barrier to adoption. Moreover, if we rely on user-provided binaries, then we have no hope  
80 of obtaining reproducible results, which in turn makes any insights far less actionable.  
81 2. It must allow for multiple revisions of the same tool to be installed simultaneously. This  
82 enables developers to easily look for performance regressions, and quantify the impact of  
83 optimizations.  
84 3. It must allow for multiple copies of the *same* version tool to be installed with different  
85 build configurations. This allows developers to look for performance regressions induced  
86 by different compiler versions/optimizations.  
87 4. It must be able to be run locally on a developer's machine. Cloud-based tools are often  
88 cumbersome to use and debug, which in turn lowers adoption.  
89 5. It must present a declarative interface for creating benchmarking environments and  
90 running benchmarks. Sandboxed builds of tools are somewhat moot if we cannot trust  
91 that a benchmark was run with the correct configuration.  
92 6. It must present performance results in a self-contained format that is easy to understand  
93 and share. Performance statistics that require large amounts of post-processing or  
94 dedicated tools to view can not be easily shared with developers, which in turn makes  
95 the data less actionable.

96 Of these criteria, the first four present the largest engineering challenge, and are tantamount  
97 to constructing a meta-build system that is able to orchestrate *other* build systems.  
98 We approached the problem by constructing a bespoke content-addressed system atop of  
99 Shake [12], which we discuss in section 5.1. The final two criteria also presented some  
100 unforeseen difficulties, which we detail in 5.2.

## 101 **5.1 The Panbench Build System**

102 As noted earlier, we strongly believe that any benchmarking system should provide infra-  
103 structure for installing multiple versions of reproducibly built software. Initially, we intended  
104 to build this infrastructure for Panbench atop of Nix [7]. This is seemingly a perfect fit;  
105 after all, Nix was designed to facilitate almost exactly this use-case. However, after further  
106 deliberation, we came to the conclusion that Nix did not quite meet our needs for the  
107 following reasons:

- 108 1. Nix does not work natively on Windows. Performance problems can be operating system  
109 specific, so ruling out an OS that has a large user base that is often overlooked in testing  
110 seems unwise<sup>3</sup>.  
111 2. Nix adds a barrier to adoption. Installing Nix is a somewhat invasive process, especially  
112 on MacOS<sup>4</sup>. We believe that it is somewhat unreasonable to ask developers to add users  
113 and modify their root directory to run a benchmarking tool, and strongly suspect that  
114 this would hamper adoption.

115 With the obvious option exhausted, we opted to create our own Nix-inspired build system  
116 based atop Shake [12]. This avoids the aforementioned problems with Nix: Shake works on  
117 Windows, and only requires potential users to install a Haskell toolchain.

---

<sup>3</sup> Currently, Panbench does not support Windows, but this is an artifact of prioritization, and not a fundamental restriction.

<sup>4</sup> The situation is even worse on x86-64 Macs, which most Nix installers simply do not support.

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

118     The details of content-addressed build systems are a deep topic unto themselves, so we  
119     will only describe the key points. Systems like Nix use an *input-addressing* scheme, wherein  
120     the results of a build are stored on disk prefixed by a hash of all build inputs. Crucially,  
121     this lets the build system know where to store the result of the build before the build is run,  
122     which avoids vicious cycles where the result of a build depends on its own hash. However,  
123     most input-addressed systems also require that the hash of the inputs *solely* determines the  
124     output. On its face, this is a reasonable ask, but taking this idea to its logical conclusion  
125     requires one to remove *any* dependency on the outer environment, which in turn forces one to  
126     re-package the entirety of the software stack all the way down to `libc`. This is an admirable  
127     goal in its own right, but is actually somewhat counterproductive for our use case: there  
128     is a very real chance that we might end up benchmarking our sub-par repackaging of some  
129     obscure C dependency four layers down the stack.

130     To avoid this cascading series of dependency issues, Panbench takes a hybrid approach,  
131     wherein builds are input-addressed, but are allowed to also depend on the external environ-  
132     ment. Additionally, the results of builds are also hashed, and stored out-of-band inside  
133     of a Shake database. This hash is used to invalidate downstream results, and also act as a  
134     fingerprint to identify if two benchmarks were created from the same binary. This enables  
135     a pay-as-you go approach to reproducibility, which we hope will result in a lower adoption  
136     cost. Moreover, we can achieve fully reproducible builds by using Nix as a meta-meta build  
137     system to compile Panbench: this is how we obtained the results presented in Section 3.

### 138 5.2 Running Benchmarks and Generating Reports

139     As noted in the introduction to this section, we believe that benchmarking tools should  
140     present a *declarative* interface for writing not just single benchmarking cases, but entire  
141     benchmarking suites and their corresponding environments. Panbench accomplishes this  
142     by *also* implementing the benchmark execution framework atop Shake. This lets us easily  
143     integrate the process of tool installation with environment setup, but introduces its own set  
144     of engineering challenges.

145     The crux of the problem is that performance tests are extremely sensitive to the current  
146     load on the system. This is largely at odds with the goals of a build system, which is to  
147     completely saturate all system resources to try to complete a build as fast as possible. This  
148     can be avoided via careful use of locks, but we are then faced with another, larger problem.  
149     Haskell is a garbage collected language, and running the GC can put a pretty heavy load on  
150     the system. Moreover, the GHC runtime system is very well engineered, and is free to run  
151     the garbage collector inside of `safe` FFI calls, and waiting for process completion is marked  
152     as `safe`.

153     To work around this, we opt to eschew existing process interaction libraries, and implement  
154     the benchmark spawning code in C<sup>5</sup>. This lets us take the rather extreme step of linking  
155     against the GHC runtime system so that we can call `rts_pause`, which pauses all other  
156     Haskell threads and GC sweeps until `rts_resume` is called.

157     Initially, we thought that this was the only concern that would arise by tightly integrating  
158     the build system with the benchmark executor. However, our initial benchmarks on Linux  
159     systems displayed some very strange behaviour, wherein the max resident set size reported  
160     by `getrusage` and `wait4` would consistently report a reading of approximately 2 gigabytes  
161     halfway through a full benchmarking suite. After some investigating, we discovered the Linux

---

<sup>5</sup> This is why Panbench does not currently support Windows.

162 preserves resource usage statistics across calls to `execve`. Consequentially, this means that  
 163 we are unable to measure any max RSS that is lower than max RSS usage of Panbench itself.  
 164 Luckily, our lowest baseline is Agda at 64 megs, and we managed to get the memory usage  
 165 of Panbench itself down to 10 megs via some careful optimization and GC tuning.

166 Currently, the statistics that Panbench gathers can then be rendered into standalone  
 167 HTML files with `vega-lite` plots, or into `TEX` files containing PGF plots. We intend to add  
 168 more visualization and statistic analysis tools as the need arises.

## 169 6 The Design of Panbench

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

176 However, the core task of Panbench is not that different from the task of a logical  
 177 framework [8]: Both exist to manipulate judgements, inference rules, and derivations:  
 178 Panbench just works with *grammatical* judgements and production rules rather than typing  
 179 judgments and inference rules. In this sense Panbench is a *grammatical* framework<sup>6</sup> rather  
 180 than a logical one.

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

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

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

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

- 192 1. A layer for defining grammars as relations.
- 193 2. A logic programming layer for synthesizing derivations.
- 194 3. A staged functional programming layer for proving “adequacy” results.

195 In this case, an adequacy result for a given language  $\mathcal{L}$  is a constructive proof that all  
 196 grammatical derivations written within the framework can be expressed within the concrete  
 197 syntax of a language  $\mathcal{L}$ . However, the computational content of such a proof essentially  
 198 amounts to a compiler written in the functional programming layer. Given that this compiler  
 199 outputs *concrete syntax*, it is implemented as a pretty-printer.

Reed: Cite something

---

<sup>6</sup> Not to be confused with *the Grammatical Framework* [13], which aims to be a more natural-language focused meta-framework for implementing and translating between grammars.

200 **6.1 Implementing The Grammatical Framework**

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

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

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

214  $\langle \text{expr} \rangle := \text{x} \mid \text{n} \mid \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \langle \text{expr} \rangle * \langle \text{expr} \rangle$   
 215  $\langle \text{stmt} \rangle := \langle \text{var} \rangle = \langle \text{expr} \rangle \mid \text{while } \langle \text{expr} \rangle \text{ do } \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmt} \rangle$

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

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

```
219 class Var expr where
220   var :: String → expr
221
222 class Lit expr where
223   lit :: Int → expr
224
225 class Add expr where
226   add :: expr → expr → expr
227
228 class Mul expr where
229   mul :: expr → expr → expr
230
231 class Assign expr stmt where
232   assign :: String → expr → stmt
233
234 class While expr stmt where
235   while :: expr → stmt → stmt
236
237 class AndThen stmt where
238   andThen :: stmt → stmt → stmt
```

240 This style of language encoding is typically known as the untyped variant of *finally*  
*tagless* [5]. However, our encoding is a slight refinement where we restrict ourselves to a single  
 241 class per production rule. Other tagless encodings often use a class per syntactic category.  
 242 This more fine-grained approach allows us to encode grammatical constructs that are only  
 243 supported by a subset of our target grammars; see section 6.2 for examples.

245 Unfortunately, the encoding above has some serious ergonomic issues. In particular,  
 246 expressions like `assign "x"` (`lit 4`) will result in an unsolved metavariable for `expr`, as there  
 247 may be multiple choices of `expr` to use when solving the `Assign ?expr stmt` constraint. Luckily,  
 248 we can resolve ambiguities of this form through judicious use of functional dependencies [10],  
 249 as demonstrated below.

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

```

250 class Assign expr stmt | stmt → expr where
251   assign :: String → expr → stmt
252
253 class While expr stmt | stmt → expr where
254   while :: expr → stmt → stmt
255

```

## 6.2 Designing The Language

Now that we've fleshed out how we are going to encode our grammatical framework into our host language, it's time to design our idealized abstract grammar. All of our target languages roughly agree on a subset of the grammar of non-binding terms: the main sources of divergence are binding forms and top-level definitions<sup>7</sup>. This is ultimately unsurprising: dependent type theories are fundamentally theories of binding and substitution, so we would expect some variation in how our target languages present the core of their underlying theories.

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

### 6.2.1 Binding Forms

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

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

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

JC: maybe a linguistic analogy would be useful? SVO vs SOV requires us to notice the syntactic categories subject, object verb that are common to all, and that explicit renderings merely differ in the order. Similarly for singular and plural as modifiers.

**Listing 3** The basic binding constructs in Panbench.

```

289 class Binder arity nm ann tm cell | cell → nm tm where
290   binder :: arity nm → ann tm → cell
291

```

<sup>7</sup> As we shall see in section 6.2.2, the syntactical divergence of top-level definitions is essentially about binders as well.

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

```

293  -- | No annotation or arity.
294  data None nm = None
295
296  -- | A single annotation or singular arity.
297  newtype Single a = Single { unSingle :: a }
298
299  -- | Multi-binders.
300  type Multi = []
301
302  -- | Infix operator for an annotated binder with a single name.
303  (.::) :: (Binder Single nm Single tm cell) => nm → tm → cell
304  nm .:: tp = binder (Single nm) (Single tp)
305
306  -- | Infix operator for an annotated binder.
307  (.*:) :: (Binder arity nm Single tm cell) => arity nm → tm → cell
308  nms .*: tp = binder nms (Single tp)

```

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

**Listing 4** The Panbench class for II-types.

```

312  class Pi cell tm | tm → cell where
313    pi :: [cell] → tm → tm
314
315

```

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

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

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

**Listing 5** Typeclasses for binding modifiers.

```

329  class Implicit cell where
330    implicit :: cell → cell
331
332
333  class SemiImplicit cell where
334    semiImplicit :: cell → cell
335

```

336 This is a case where the granular tagless encoding shines. Both Lean 4 and Rocq have a  
 337 form of semi-implicits<sup>8</sup>, whereas Idris 2 and Agda have no such notion. Decomposing the  
 338 language of binding modifiers into granular pieces lets us write benchmarks that explicitly  
 339 require support for features like semi-implicits. Had we used a monolithic class that encodes

---

<sup>8</sup> 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.

340 the entire language of modifiers, we would have to resort to runtime errors (or, even worse,  
341 dubious attempts at translation).

### 342 6.2.2 Top-Level Definitions

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

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

```
348 349 private instance abstract @irr @mixed foo bar : Nat → _
```

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

353 Unfortunately, this insight does not offer a complete solution. Notably, our four target  
354 grammar differ significantly in how their treatment of type signatures. prioritize dependent  
355 pattern matching (e.g. Agda, Idris 2) typically opt to have standalone type signatures: this  
356 allow for top-level pattern matches, which in turn makes it much easier to infer motives[11].  
357 Conversely, languages oriented around tactics (e.g. Lean 4, Rocq) typically opt for in-line  
358 type signatures and pattern-matching expressions. This appears to be largely independent of  
359 Miranda/ML syntax split: Lean 4 borrows large parts of its syntax from Haskell, yet still  
360 opts for in-line signatures.

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

**■ Listing 7** A definition with mismatched names.

```
366 367 id : (A : Type) → A → A
368 id B x = x
```

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

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

**■ Listing 8** Two variants of the identity function.

```
378
379 id : {A : Type} → A → A
380 id x = x
381
382 id' : {A : Type} → A → A
383 id' {A} x = x
```

## 23:10 Panbench: A Comparative Benchmarking Tool for Dependently-Typed Languages

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

391 We have encoded this decision in our idealized grammar by introducing a notion of a  
392 “left-hand-side” of a definition, which consists of a collection of names to be defined, and  
393 a scope to define them under. This means that we view definitions like Listing 8 not as  
394 functions `id : (A : Type) → A → A` but rather as *bindings* `A : Type, x : A ⊢ id : A` in  
395 non-empty contexts. This shift in perspective has the added benefit of making the interface  
396 to other forms of parameterised definitions entirely uniform; for instance, a parameterised  
397 record is simply just a record with a non-empty left-hand side.

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

### Listing 9 Definitions and left-hand sides.

```
400 class Definition lhs tm defn | defn → lhs tm where
401   (. =) :: lhs → tm → defn
402
403 class DataDefinition lhs ctor defn | defn → lhs ctor where
404   data_ :: lhs → [ctor] → defn
405
406 class RecordDefinition lhs nm fld defn | defn → lhs nm fld where
407   record_ :: lhs → name → [fld] → defn
```

JC: should  
have a closing  
sentence

## 7 Conclusion

### References

- 413 1 February 2026. URL: <https://github.com/GrammaticalFramework/informath>.
- 414 2 Agda Developers. Agda. URL: <https://agda.readthedocs.io/>.
- 415 3 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.
- 416 4 Edwin Brady. Idris 2: Quantitative Type Theory in Practice. pages 32460–33960, 2021. doi:10.4230/LIPICS.ECOOP.2021.9.
- 417 5 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.
- 418 6 Leonardo de Moura and Sebastian Ullrich. The lean 4 theorem prover and programming language. In *Automated Deduction – CADE 28: 28th International Conference on Automated Deduction, Virtual Event, July 12–15, 2021, Proceedings*, page 625–635, Berlin, Heidelberg, 2021. Springer-Verlag. doi:10.1007/978-3-030-79876-5\_37.
- 419 7 Eelco Dolstra and The Nix contributors. Nix. URL: <https://github.com/NixOS/nix>.
- 420 8 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.
- 421 9 Hugo Herbelin, Pierre-Marie Pédrot, coqbot, Gaëtan Gilbert, Maxime Dénès, letouzey, Emilio Jesús Gallego Arias, Matthieu Sozeau, Théo Zimmermann, Enrico Tassi, Jean-Christophe

- 434 Filliatre, Pierre Roux, Guillaume Melquiond, Arnaud Spiwack, Jason Gross, barras, Pierre  
435 Boutillier, Vincent Laporte, Jim Fehrle, Stéphane Glondu, Yves Bertot, Jasper Hugunin,  
436 Clément Pit-Claudel, Ali Caglayan, forestjulien, Pierre Courtieu, Frédéric Besson, Pierre  
437 Rousselin, MSoegtropIMC, and Yann Leray. Rocq, February 2026. doi:10.5281/zenodo.  
438 1003420.
- 439 10 Mark P. Jones. Type classes with functional dependencies. In Gert Smolka, editor, *Programming  
440 Languages and Systems*, page 230–244, Berlin, Heidelberg, 2000. Springer. doi:10.1007/  
441 3-540-46425-5\_15.
- 442 11 Conor McBride and James McKinna. The view from the left. *Journal of Functional Programming*,  
443 14(1):69–111, January 2004. doi:10.1017/S0956796803004829.
- 444 12 Neil Mitchell. Shake before building: replacing make with haskell. *SIGPLAN Not.*, 47(9):55–66,  
445 2012. doi:10.1145/2398856.2364538.
- 446 13 Aarne Ranta. *Grammatical framework: programming with multilingual grammars*. CSLI  
447 studies in computational linguistics. CSLI Publications, Center for the Study of Language and  
448 Information, Stanford, Calif, 2011.