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

**Anonymous author**
Anonymous affiliation

**Anonymous author**
Anonymous affiliation

## Abstract

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

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

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

## 1 Introduction

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

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

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

1. An extensible embedded Haskell DSL that encodes the concrete syntax of a typical dependently typed language.

2. A series of compilers for that DSL to Agda, Idris 2, Lean 4, and Rocq.

3. A benchmarking harness that can perform sandboxed builds of multiple versions of Agda, Idris 2, Lean 4, and Rocq.

4. An incremental build system that can produce benchmarking reports as static HTML files or PGF plots[1].

Note that we uniformly report all results in alphabetical order of the system's names.

Reed: Citations here

## 2 Methodology

To perform reliable benchmarking, we need tooling and we need to design a test suite. For reproducibility, we also need to document the actual experimental setup. We will describe the tooling in detail in later sections, for now we focus on the design of the test suite and the experimental setup.

## 2.1 Designing tests

Let us assume that we possess the following tools:

- a single language of tests,
- a means to translate these tests to each of our four languages,
- a reproducible test harness,
- a stable installation of the four languages.

How would we design actual tests for common features?

As we said before, our aim here is to test "basic engineering". For example, how does each system handle giant files (e.g.: one million empty lines), large names (thousands of characters long) in each syntactic category, large numbers of fields in records, constructors in algebraic types, large numbers of parameters or indices, long dependency chains, and so on. We divide our tests into the following categories:

1. basic syntactic capabilities
2. long names in simple structures
3. large simple structures (data, record)
4. handling of nesting and dependency
5. conversion, postulates

where we vary the "size" of each. Table 1 gives more details of the available tests.

At a higher level, we asked ourselves the question: for each aspect of a dependently typed language (lexical structure, grammatical structure, fundamental features of dependent languages), what could be made to "scale"? The only features of interest remained the ones that were clearly features of all four systems.

For lack of space, we cannot show samples of all the tests, but we show just a few that should be *illustrative* of the rest[2]. The source of all tests can be found in the accompanying material. All tests below are from the snapshot of our "golden" test suite, uniformly run with a size parameter of 5.

---

[1] All plots in this paper were produced directly by Panbench.
[2] Very minor edits made to fit the page

| Category | Details | Category | Details |
|---|---|---|---|
| Syntax | Newlines | Datatypes | Parameters |
| | Parentheses | | Indices |
| Names | Datatype | | Constructors |
| | Data constructor | | Params + Indices |
| | Definition | Records | Fields |
| | Definition lhs | | Parameters |
| | Definition rhs | Dependency | Record Fields |
| | Lambda | | Datatypes Parameters |
| | Pi | | Definitions chains |
| | Record constructor | | Record chains |
| | Record field | Nesting | Id chain |
| Binders | Lambda | | Id chain $\lambda$ |
| | Implicits | | Let |
| Misc | Postulates | | Let addition |
| Conversion | Addition | | Let functions |

**Figure 1** Our tests

**Listing 1** "Nesting: Id chain (Agda)"

```
module IdChain where

f : {a : Set} (x : a) → a
f x = x

test : {a : Set} → a → a
test = f f f f f f
```

**Listing 2** "Nesting: Let add (Lean)"

```
def n : Nat :=
  let x0 := 1
  let x1 := x0 + x0
  let x2 := x1 + x1
  let x3 := x2 + x2
  let x4 := x3 + x3
  x4
```

**Listing 3** "Binders: Lambda (Idris 2)"

```
module Main

const5 : {A : Type} ->
         {B0, B1, B2, B3, B4, B5 : Type} ->
         A -> B0 -> B1 -> B2 -> B3 -> B4 -> B5 -> A
const5 = \a, b0, b1, b2, b3, b4, b5 => a
```

**Listing 4** "Dependency: Record chains (Rocq)"

```
Module SequentialDependentRecords.
```

```
96   Record Dummy1 : Type := Const1 { f1 : nat }.
97   Record Dummy2 : Type := Const2 { f2 : Dummy1 }.
98   Record Dummy3 : Type := Const3 { f3 : Dummy2 }.
99   Record Dummy4 : Type := Const4 { f4 : Dummy3 }.
100  Record Dummy5 : Type := Const5 { f5 : Dummy4 }.
101
102  Definition test : Dummy5 :=
103    (Const5 (Const4 (Const3 (Const2 (Const1 10)))))).
104
105  End SequentialDependentRecords.
106
```

## 2.2   Experimental setup

All tests are run on a dedicated (older) desktop machine running on NixOS 25.11. The CPU is CPU is a an Intel i7-2600K with 1MB L2 cache, 8MB L3 cache, running at 3.4 Ghz, equipped with 24 Gigs of DDR3 memory. This box has no SSD drive, but this should only affect system time.

When tests are run, no other (human) activity happens on the machine.

All tests are run with a time limit of 60 seconds. We also tried to use memory limits but, alas, Linux no longer reliably supports these. Otherwise, we ran all four systems in their default configuration.

> JC: Reed, please fill-in X and Y below.

When we report times, we report X. We also report "max RSS" (maximum resident set size) as given by Y.

## 3   Results

> JC: I say 32 tests here, table has 30. Which one is wrong?

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

| System | User Time (s) | System Time (s) | Max RSS (MB) |
|--------|---------------|-----------------|--------------|
| Agda   | 0.02 (0.002)  | 0.01 (0.001)    | 64 (0.1)     |
| Idris 2 | 0.58 (0.007) | 0.10 (0.007)    | 248 (0.1)    |
| Lean 4 | 0.14 (0.005)  | 0.04 (0.005)    | 307 (0.6)    |
| Rocq   | 0.05 (0.004)  | 0.03 (0.003)    | 95 (0.05)    |

**Figure 2** Start-up time and memory, mean with standard deviation in parentheses.

## 4   Discussion

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

---

[3] Nor can we have appendices!

**4.0.0.1 General**

**4.0.0.2 Agda**

**4.0.0.3 Idris 2**

**4.0.0.4 Lean 4**

**4.0.0.5 Rocq**

## 5 Infrastructure

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

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

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

### 5.1 The Panbench Build System

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

1. Nix does not work natively on Windows. Performance problems can be operating system specific, so ruling out an OS that has a large user base that is often overlooked in testing seems unwise[4].

2. Nix adds a barrier to adoption. Installing Nix is a somewhat invasive process, especially on MacOS[5]. We believe that it is somewhat unreasonable to ask developers to add users and modify their root directory to run a benchmarking tool, and strongly suspect that this would hamper adoption.

With the obvious option exhausted, we opted to create our own Nix-inspired build system based atop Shake [8]. This avoids the aforementioned problems with Nix: Shake works on Windows, and only requires potential users to install a Haskell toolchain.

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

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

## 5.2 Running Benchmarks and Generating Reports

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

The crux of the problem is that performance tests are extremely sensitive to the current load on the system. This is largely at odds with the goals of a build system, which is to completely saturate all system resources to try to complete a build as fast as possible. This can be avoided via careful use of locks, but we are then faced with another, larger problem. Haskell is a garbage collected language, and running the GC can put a pretty heavy load on

---

[4] Currently, Panbench does not support Windows, but this is an artifact of prioritization, and not a fundamental restriction.

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

the system. Moreover, the GHC runtime system is very well engineered, and is free to run the garbage collector inside of `safe` FFI calls, and waiting for process completion is marked as `safe`.

To work around this, we opt to eschew existing process interaction libraries, and implement the benchmark spawning code in C[6]. This lets us take the rather extreme step of linking against the GHC runtime system so that we can call `rts_pause`, which pauses all other Haskell threads and GC sweeps until `rts_resume` is called.

Initially, we thought that this was the only concern that would arise by tightly integrating the build system with the benchmark executor. However, our initial benchmarks on Linux systems displayed some very strange behaviour, wherein the max resident set size reported by `getrusage` and `wait4` would consistently report a reading of approximately 2 gigabytes halfway through a full benchmarking suite. After some investigating, we discovered the Linux preserves resource usage statistics across calls to `execve`. Consequentially, this means that we are unable to measure any max RSS that is lower than max RSS usage of Panbench itself. Luckily, our lowest baseline is Agda at 64 megs, and we managed to get the memory usage of Panbench itself down to 10 megs via some careful optimization and GC tuning.

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

## 6 The Design of Panbench

At its core, Panbench is a tool for producing grammatically well-formed concrete syntax across multiple different languages. Crucially, Panbench does *not* require that the syntax produced is well-typed or even well-scoped: if it did, then it would be impossible to benchmark 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.

Reed: Cite something

However, the core task of Panbench is not that different from the task of a logical framework [5]: Both exist to manipulate judgements, inference rules, and derivations: Panbench just works with *grammatical* judgements and production rules rather than typing judgments and inference rules. In this sense Panbench is a *grammatical* framework[7] rather than a logical one.

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

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

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

---

[6] This is why Panbench does not currently support Windows.
[7] 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.

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

**1.** A layer for defining grammars as relations.
**2.** A logic programming layer for synthesizing derivations.
**3.** A staged functional programming layer for proving "adequacy" results.

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

## 6.1    Implementing The Grammatical Framework

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

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

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

$\langle expr \rangle := \mathrm{x} \mid \mathrm{n} \mid \langle expr \rangle \ `\texttt{+}` \ \langle expr \rangle \mid \langle expr \rangle \ `\texttt{*}` \ \langle expr \rangle$

$\langle stmt \rangle := \langle var \rangle \ `\texttt{=}` \ \langle expr \rangle \mid `\texttt{while}` \ \langle expr \rangle \ `\texttt{do}` \ \langle stmt \rangle \mid \langle stmt \rangle \ `\texttt{;}` \ \langle stmt \rangle$

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

**Listing 5** An example tagless encoding.

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

class Lit expr where
  lit :: Int → expr

class Add expr where
  add :: expr → expr → expr

class Mul expr where
  mul :: expr → expr → expr

class Assign expr stmt where
  assign :: String → expr → stmt

class While expr stmt where
  while :: expr → stmt → stmt
```

```
299  class AndThen stmt where
300    andThen :: stmt → stmt → stmt
301
```

This style of language encoding is typically known as the untyped variant of *finally tagless* [3]. However, our encoding is a slight refinement where we restrict ourselves to a single class per production rule. Other tagless encodings often use a class per syntactic category. This more fine-grained approach allows us to encode grammatical constructs that are only supported by a subset of our target grammars; see section 6.2 for examples.

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

▪ **Listing 6** A tagless encoding with functional dependencies.

```
312
313  class Assign expr stmt | stmt → expr where
314    assign :: String → expr → stmt
315
316  class While expr stmt | stmt → expr where
317    while :: expr → stmt → stmt
318
```

## 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[8]. 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 $(x\ y\ z : A) \to B$. 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

---

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

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 singu-... as modifiers.

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.

▨ **Listing 7** The basic binding constructs in Panbench.

```
class Binder arity nm ann tm cell | cell → nm tm where
  binder :: arity nm → ann tm → cell

-- | No annotation or arity.
data None nm = None

-- | A single annotation or singular arity.
newtype Single a = Single { unSingle :: a }

-- | Multi-binders.
type Multi = []

-- | Infix operator for an annotated binder with a single name.
(.:) :: (Binder Single nm Single tm cell) => nm → tm → cell
nm .: tp = binder (Single nm) (Single tp)

-- | Infix operator for an annotated binder.
(.:*) :: (Binder arity nm Single tm cell) => arity nm → tm → cell
nms .:* tp = binder nms (Single tp)
```

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

▨ **Listing 8** The Panbench class for Π-types.

```
class Pi cell tm | tm → cell where
  pi :: [cell] → tm → tm
```

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

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

The language of binding modifiers is implemented as the following set of Haskell typeclasses.

**Listing 9** Typeclasses for binding modifiers.

```
class Implicit cell where
   implicit :: cell → cell

class SemiImplicit cell where
   semiImplicit :: cell → cell
```

This is a case where the granular tagless encoding shines. Both Lean 4 and Rocq have a form of semi-implicits[9], whereas Idris 2 and Agda have no such notion. Decomposing the language of binding modifiers into granular pieces lets us write benchmarks that explicitly require support for features like semi-implicits. Had we used a monolithic class that encodes the entire language of modifiers, we would have to resort to runtime errors (or, even worse, dubious attempts at translation).

### 6.2.2 Top-Level Definitions

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

**Listing 10** A complicated Agda signature.

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

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

Unfortunately, this insight does not offer a complete solution. Notably, our four target grammar differ significantly in how their treatment of type signatures. prioritize dependent pattern matching (e.g. Agda, Idris 2) typically opt to have standalone type signatures: this allow for top-level pattern matches, which in turn makes it much easier to infer motives[7]. Conversely, languages oriented around tactics (e.g. Lean 4, Rocq) typically opt for in-line type signatures and pattern-matching expressions. This appears to be largely independent of Miranda/ML syntax split: Lean 4 borrows large parts of its syntax from Haskell, yet still opts for in-line signatures.

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

**Listing 11** A definition with mismatched names.

```
   id : (A : Type) → A → A
   id B x = x
```

In particular, note that we have bound the first argument to a different name. Translating this to a corresponding Rocq declaration then forces us to choose to use either the name from

---

[9] 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.

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 12** Two variants of the identity function.

```
id : {A : Type} → A → A
id x = x

id' : {A : Type} → A → A
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 Listing 12 not as functions id : $(A : \text{Type}) \to A \to A$ but rather as *bindings* $A : \text{Type}, x : A \vdash \text{id} : A$ 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.

**Listing 13** Definitions and left-hand sides.

```
class Definition lhs tm defn | defn → lhs tm where
  (.=) :: lhs → tm → defn

class DataDefinition lhs ctor defn | defn → lhs ctor where
  data_ :: lhs → [ctor] → defn

class RecordDefinition lhs nm fld defn | defn → lhs nm fld where
  record_ :: lhs → name → [fld] → defn
```

JC: should have a closing sentence

## 7    Conclusion

### References

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

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.

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`.

4   Eelco Dolstra and The Nix contributors. Nix. URL: `https://github.com/NixOS/nix`.

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`.

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`.

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`.

8   Neil Mitchell. Shake before building: replacing make with haskell. *SIGPLAN Not.*, 47(9):55–66, 2012. `doi:10.1145/2398856.2364538`.

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.