

The Backpack algorithm

July 7, 2015

This document describes the Backpack shaping and typechecking passes, as we intend to implement it.

1 Changelog

April 28, 2015 A signature declaration no longer provides a signature in the technical shaping sense; the motivation for this change is explained in **In-scope signatures are not provisions**. The simplest consequence of this is that all requirements are importable (Derek has stated that he doesn't think this will be too much of a problem in practice); it is also possible to extend shape with a `signatures` field, although some work has to be done specifying coherence conditions between `signatures` and `requirements`.

2 Front-end syntax

p, q, r	Package names	
m, n	Module names	
Packages		
pkg	<code>::= package p [$provreq$] where $\{d_1; \dots; d_n\}$</code>	
Declarations		
d	<code>::= module m [$exports$] where $body$</code> <code>signature m [$exports$] where $body$</code> <code>include p [$provreq$]</code>	
Provides/requires specification		
$provreq$	<code>::= (rns) [<code>requires</code>(rns)]</code>	
rns	<code>::= rn_0, \dots, rn_n[,]</code>	Renamings
rn	<code>::= m as n</code>	Renaming
Haskell code		
$exports$	A Haskell module export list	
$body$	A Haskell module body	

Figure 1: Syntax of Backpack

The syntax of Backpack is given in Figure 1. See the “Backpack manual” for more explanation about the syntax. It is slightly simplified here by removing any constructs which are easily implemented as syntactic sugar (e.g., a bare m in a renaming is simply `m as m` .)

3 Shaping

```
Shape ::= provides: m -> Module { Name, ... }; ...
       requires: m ->      { Name, ... }; ...
PkgKey ::= p(m -> Module, ... )
Module ::= PkgKey:m
Name    ::= Module.OccName
OccName Unqualified name in a namespace
```

Figure 2: Semantic entities in Backpack

Shaping computes a *Shape*, whose form is described in Figure 2. A shape describes what modules a package implements and exports (the *provides*) and what signatures a package needs to have filled in (the *requires*). Both provisions and requires can be imported after a package is included.

We incrementally build a shape by starting with an empty shape context and adding to it as follows:

1. Calculate the shape of a declaration, with respect to the current shape context. (e.g., by renaming a module/signature, or using the shape from an included package.)
2. Merge this shape into the shape context.

The final shape context is the shape of the package as a whole. Optionally, we can also compute the renamed syntax trees of modules and signatures.

In the description below, we'll assume **THIS** is the package key of the package being processed.

3.1 module M

A module declaration provides a module **THIS:M** at module name **M**. It has the shape:

```
provides: M -> THIS:M { exports of renamed M under THIS:M }
requires: (nothing)
```

Example:

```
module A(T) where
  data T = T

-- provides: A -> THIS:A { THIS:A.T }
-- requires: (nothing)
```

OccName is implied by *Name*. In Haskell, the following is not valid syntax:

```
import A (foobar as baz)
```

In particular, a *Name* which is in scope will always have the same *OccName* (even if it may be qualified.) You might imagine relaxing this restriction so that declarations can be used under different *OccNames*; in such a world, we need a different definition of shape:

```
Shape ::=
  provided: ModName -> { OccName -> Name }
  required: ModName -> { OccName -> Name }
```

Presently, however, such an *OccName* annotation would be redundant: it can be inferred from the *Name*.

Holes of a package are a mapping, not a set. Why can't the *PkgKey* just record a set of *Modules*, e.g. *PkgKey ::= SrcPkgKey { Module }*? Consider:

```
package p (A) requires (H1, H2) where
  signature H1(T) where
    data T
  signature H2(T) where
    data T
  module A(A(..)) where
    import qualified H1
    import qualified H2
    data A = A H1.T H2.T

package q (A12, A21) where
  module I1(T) where
    data T = T Int
  module I2(T) where
    data T = T Bool
  include p (A as A12) requires (H1 as I1, H2 as I2)
  include p (A as A21) requires (H1 as I2, H2 as I1)
```

With a mapping, the first instance of *p* has key *p*(*H1* -> *q*() : *I1*, *H2* -> *q*() : *I2*) while the second instance has key *p*(*H1* -> *q*() : *I2*, *H2* -> *q*() : *I1*); with a set, both would have the key *p*(*q*() : *I1*, *q*() : *I2*).

Signatures can require a specific entity. With requirements like *A* -> { *HOLE:A.T*, *HOLE:A.foo* }, why not specify it as *A* -> { *T*, *foo* }, e.g., *required: { ModName -> { OccName } }*? Consider:

```
package p () requires (A, B) where
  signature A(T) where
    data T
  signature B(T) where
    import T
```

The requirements of this package specify that *A.T* = *B.T*; this can be expressed with *Names* as

```
A -> { HOLE:A.T }
B -> { HOLE:A.T }
```

But, without *Names*, the sharing constraint is impossible: *A* -> { *T* }; *B* -> { *T* }. (NB: *A* and *B* don't have to be implemented with the same module.)

The *Name* of a value is used to avoid ambiguous identifier errors. We state that two types are equal when their *Names* are the same; however, for values, it is less clear why we care. But consider this example:

```
package p (A) requires (H1, H2) where
  signature H1(x) where
    x :: Int
  signature H2(x) where
    import H1(x)
  module A(y) where
    import H1
    import H2
    y = x
```

The reference to `x` in `A` is unambiguous, because it is known that `x` from `H1` and `x` from `H2` are the same (have the same *Name*.) If they were not the same, it would be ambiguous and should cause an error. Knowing the *Name* of a value distinguishes between these two cases.

Holes are linear Requirements do not record what *Module* represents the identity of a requirement, which means that it's not possible to assert that hole `A` and hole `B` should be implemented with the same module, as might occur with aliasing:

```
signature A where
signature B where
alias A = B
```

The benefit of this restriction is that when a requirement is filled, it is obvious that this is the only requirement that is filled: you won't magically cause some other requirements to be filled. The downside is it's not possible to write a package which looks for an interface it is looking for in one of n names, accepting any name as an acceptable linkage. If aliasing was allowed, we'd need a separate physical shaping context, to make sure multiple mentions of the same hole were consistent.

3.2 signature M

A signature declaration creates a requirement at module name `M`. It has the shape:

```
provides: (nothing)
requires: M -> { exports of renamed M under HOLE:M }
```

Example:

```
signature H(T) where
  data T

-- provides: H -> (nothing)
-- requires: H -> { HOLE:H.T }
```

In-scope signatures are not provisions. We enforce the invariant that a provision is always (syntactically) a **module** and a requirement is always a **signature**. This means that if you have a requirement and a provision of the same name, the requirement can *always* be filled with the provision. Without this invariant, it's not clear if a provision will actually fill a signature. Consider this example, where a signature is required and exposed:

```
package a-sigs (A) requires (A) where -- ***
  signature A where
    data T

package a-user (B) requires (A) where
  signature A where
    data T
    x :: T
  module B where
    ...

package p where
  include a-sigs
  include a-user
```

When we consider merging in the shape of `a-user`, does the `A` provided by `a-sigs` fill in the `A` requirement in `a-user`? It *should not*, since `a-sigs` does not actually provide enough declarations to satisfy `a-user`'s requirement: the intended semantics *merges* the requirements of `a-sigs` and `a-user`.

```
package a-sigs (M as A) requires (H as A) where
  signature H(T) where
    data T
  module M(T) where
    import H(T)
```

We rightly should error, since the provision is a module. And in this situation:

```
package a-sigs (H as A) requires (H) where
  signature H(T) where
    data T
```

The requirements should be merged, but should the merged requirement be under the name `H` or `A`? It may still be possible to use the `(A) requires (A)` syntax to indicate exposed signatures, but this would be a mere syntactic alternative to `() requires (exposed A)`.

3.3 include pkg (X) requires (Y)

We merge with the transformed shape of package `pkg`, where this shape is transformed by:

- Renaming and thinning the provisions according to (X)
- Renaming requirements according to (Y) (requirements cannot be thinned, so non-mentioned requirements are implicitly passed through.) For each renamed requirement from Y to Y', substitute `HOLE:Y` with `HOLE:Y'` in the *Modules* and *Names* of the provides and requires.

If there are no thinnings/renamings, you just merge the shape unchanged! Here is an example:

```
package p (M) requires (H) where
  signature H where
    data T
  module M where
    import H
    data S = S T

package q (A) where
  module X where
    data T = T
  include p (M as A) requires (H as X)
```

The shape of package `p` is:

```
requires: M -> { p(H -> HOLE:H):M.S }
provides: H -> { HOLE:H.T }
```

Thus, when we process the `include` in package `q`, we make the following two changes: we rename the provisions, and we rename the requirements, substituting `HOLEs`. The resulting shape to be merged in is:

```
provides: A -> { p(H -> HOLE:X):M.S }
requires: X -> { HOLE:X.T }
```

After merging this in, the final shape of `q` is:

```
provides: X -> { q():X.T }           -- from shaping 'module X'
          A -> { p(H -> q():X):M.S }
requires: (nothing)                 -- discharged by provided X
```

3.4 Merging

The shapes we've given for individual declarations have been quite simple. Merging combines two shapes, filling requirements with implementations, unifying *Names*, and unioning requirements; it is the most complicated part of the shaping process.

The best way to think about merging is that we take two packages with inputs (requirements) and outputs (provisions) and “wiring” them up so that outputs feed into inputs. In the absence of mutual recursion, this wiring process is *directed*: the provisions of the first package feed into the requirements of the second package, but never vice versa. (With mutual recursion, things can go in the opposite direction as well.)

Suppose we are merging shape p with shape q (e.g., $p; q$). Merging proceeds as follows:

1. *Fill every requirement of q with provided modules from p .* For each requirement M of q that is provided by p (in particular, all of its required **Names** are provided), substitute each *Module* occurrence of `HOLE:M` with the provided $p(M)$, unify the names, and remove the requirement from q . If the names of the provision are not a superset of the required names, error.
2. If mutual recursion is supported, *fill every requirement of p with provided modules from q .*
3. *Merge leftover requirements.* For each requirement M of q that is not provided by p but required by p , unify the names, and union them together to form the new requirement. (It's not necessary to substitute *Modules*, since they are guaranteed to be the same.)
4. *Add provisions of q .* Union the provisions of p and q , erroring if there is a duplicate that doesn't have the same identity.

To unify two sets of names, find each pair of names with matching *OccNames* n and m and do the following:

1. If both are from holes, pick a canonical representative m and substitute n with m .
2. If one n is from a hole, substitute n with m .
3. Otherwise, error if the names are not the same.

It is important to note that substitutions on *Modules* and substitutions on *Names* are disjoint: a substitution from `HOLE:A` to `HOLE:B` does *not* substitute inside the name `HOLE:A.T`.

Since merging is the most complicated step of shaping, here are a big pile of examples of it in action.

3.4.1 A simple example

In the following set of packages:

```
package p(M) requires (A) where
  signature A(T) where
    data T
  module M(T, S) where
    import A(T)
    data S = S T

package q where
  module A where
    data T = T
  include p
```

When we `include p`, we need to merge the partial shape of `q` (with just provides `A`) with the shape of `p`. Here is each step of the merging process:

shape 1	shape 2

(initial shapes)	
provides: A -> THIS:A { q():A.T }	M -> p(A -> HOLE:A) { HOLE:A.T, p(A -> HOLE:A).S }
requires: (nothing)	A -> { HOLE:A.T }
(after filling requirements)	
provides: A -> THIS:A { q():A.T }	M -> p(A -> THIS:A) { q():A.T, p(A -> THIS:A).S }
requires: (nothing)	(nothing)
(after adding provides)	
provides: A -> THIS:A { q():A.T }	
	M -> p(A -> THIS:A) { q():A.T, p(A -> THIS:A).S }
requires: (nothing)	

Notice that we substituted HOLE:A with THIS:A, but HOLE:A.T with q():A.T.

3.4.2 Requirements merging can affect provisions

When a merge results in a substitution, we substitute over both requirements and provisions:

```
signature H(T) where
  data T
module A(T) where
  import H(T)
module B(T) where
  data T = T

-- provides: A -> THIS:A { HOLE:H.T }
--           B -> THIS:B { THIS:B.T }
-- requires: H ->      { HOLE:H.T }

signature H(T, f) where
  import B(T)
  f :: a -> a

-- provides: A -> THIS:A { THIS:B.T }           -- UPDATED
--           B -> THIS:B { THIS:B.T }
-- requires: H ->      { THIS:B.T, HOLE:H.f } -- UPDATED
```

3.4.3 Sharing constraints

Suppose you have two signature which both independently define a type, and you would like to assert that these two types are the same. In the ML world, such a constraint is known as a sharing constraint. Sharing constraints can be encoded in Backpacks via clever use of reexports; they are also an instructive example for signature merging.

```
signature A(T) where
  data T
signature B(T) where
  data T

-- requires: A -> { HOLE:A.T }
```



```

      B -> { HOLE:B.T }

-- the sharing constraint!
signature A(T) where
  import B(T)
-- (shape to merge)
-- requires: A -> { HOLE:B.T }

-- (after merge)
-- requires: A -> { HOLE:A.T }
--           B -> { HOLE:A.T }

```

I'm pretty sure any choice of *Name* is OK, since the subsequent substitution will make it alpha-equivalent.

3.5 Export declarations

If an explicit export declaration is given, the final shape is the computed shape, minus any provisions not mentioned in the list, with the appropriate renaming applied to provisions and requirements. (Requirements are implicitly passed through if they are not named.) If no explicit export declaration is given, the final shape is the computed shape, including only provisions which were defined in the declarations of the package.

3.6 Package key

What is THIS? It is the package name, plus for every requirement M, a mapping M -> HOLE:M. Annoyingly, you don't know the full set of requirements until the end of shaping, so you don't know the package key ahead of time; however, it can be substituted at the end easily.

Signature visibility, and defaulting The simplest formulation of requirements is to have them always be visible. Signature visibility could be controlled by associating every requirement with a flag indicating if it is importable or not: a signature declaration sets a requirement to be visible, and an explicit export list can specify if a requirement is to be visible or not.

When an export list is absent, we have to pick a default visibility for a signature. If we use the same behavior as with modules, a strange situation can occur:

```
package p where -- S is visible
  signature S where
    x :: True

package q where -- use defaulting
  include p
  signature S where
    y :: True
  module M where
    import S
    z = x && y      -- OK

package r where
  include q
  module N where
    import S
    z = y          -- OK
    z = x          -- ???
```

Absent the second signature declaration in `q`, `S.x` clearly should not be visible in `N`. However, what ought to occur when this signature declaration is added? One interpretation is to say that only some (but not all) declarations are provided (`S.x` remains invisible); another interpretation is that adding `S` is enough to treat the signature as “in-line”, and all declarations are now provided (`S.x` is visible).

The latter interpretation avoids having to keep track of providedness per declarations, and means that you can always express defaulting behavior by writing an explicit provides declaration on the package. However, it has the odd behavior of making empty signatures semantically meaningful:

```
package q where
  include p
  signature S where
```

4 Type constructor exports

In the previous section, we described the *Names* of a module as a flat namespace; but actually, there is one level of hierarchy associated with type-constructors. The type:

```
data A = B { foo :: Int }
```

brings three *OccNames* into scope, *A*, *B* and *foo*, but the constructors and record selectors are considered *children* of *A*: in an import list, they can be implicitly brought into scope with *A(..)*, or individually brought into scope with *foo* or *pattern B* (using the new *PatternSynonyms* extension). Symmetrically, a module may export only *some* of the constructors/selectors of a type; it may not even export the type itself!

We *absolutely* need this information to rename a module or signature, which means that there is a little bit of extra information we have to collect when shaping. What is this information? If we take GHC's internal representation at face value, we have the more complex semantic representation seen in Figure 3:

```
Shape ::= provides: m -> Module { AvailInfo, ... }; ...
        requires: m ->      { AvailInfo, ... }; ...
AvailInfo ::= Name                                     Plain identifiers
           | Name { Name0, ..., Namen }             Type constructors
```

Figure 3: Enriched semantic entities in Backpack

For type constructors, the outer *Name* identifies the *parent* identifier, which may not necessarily be in scope (define this to be the *availName*); the inner list consists of the children identifiers that are actually in scope. If a wildcard is written, all of the child identifiers are brought into scope. In the following examples, we've ensured that types and constructors are unambiguous, although in Haskell proper they live in separate namespaces; we've also elided the *THIS* package key from the identifiers.

```
module M(A(..)) where
  data A = B { foo :: Int }
-- M.A{ M.A, M.B, M.foo }

module N(A) where
  data A = B { foo :: Int }
-- N.A{ N.A }

module O(foo) where
  data A = B { foo :: Int }
-- O.A{ O.foo }

module A where
  data T = S { bar :: Int }
module B where
  data T = S { baz :: Bool }
module C(bar, baz) where
  import A(bar)
  import B(baz)
-- A.T{ A.bar }, B.T{ B.baz }
-- NB: it would be illegal for the type constructors
-- A.T and B.T to be both exported from C!
```

Previously, we stated that we simply merged *Names* based on their *OccNames*. We now must consider what it means to merge *AvailInfos*.

4.1 Algorithm

Our merging algorithm takes two sets of *AvailInfos* and merges them into one set. In the degenerate case where every *AvailInfo* is a *Name*, this algorithm operates the same as the original algorithm. Merging proceeds in two steps: unification and then simple union.

Unification proceeds as follows: for each pair of *Names* with matching *OccNames*, unify the names. For each pair of *Name* $\{ Name_0, \dots, Name_n \}$, where there exists some pair of child names with matching *OccNames*, unify the parent *Names*. (A single *AvailInfo* may participate in multiple such pairs.) A simple identifier and a type constructor *AvailInfo* with overlapping in-scope names fails to unify. After unification, the simple union combines entries with matching *availNames* (parent name in the case of a type constructor), recursively unioning the child names of type constructor *AvailInfos*.

Unification of *Names* results in a substitution, and a *Name* substitution on *AvailInfo* is a little unconventional. Specifically, substitution on *Name* $\{ Name_0, \dots, Name_n \}$ proceeds specially: a substitution from *Name* to *Name'* induces a substitution from *Module* to *Module'* (as the *OccNames* of the *Names* are guaranteed to be equal), so for each child *Name_i*, perform the *Module* substitution. So for example, the substitution `HOLE:A.T` to `THIS:A.T` takes the *AvailInfo* `HOLE:A.T { HOLE:A.B, HOLE:A.foo }` to `THIS:A.T { THIS:A.B, THIS:A.foo }`. In particular, substitution on children *Names* is *only* carried out by substituting on the outer name; we will never directly substitute children.

4.2 Examples

Unfortunately, there are a number of tricky scenarios:

Merging when type constructors are not in scope

```
signature A1(foo) where
  data A = A { foo :: Int, bar :: Bool }

signature A2(bar) where
  data A = A { foo :: Int, bar :: Bool }
```

If we merge `A1` and `A2`, are we supposed to conclude that the types `A1.A` and `A2.A` (not in scope!) are the same? The answer is no! Consider these implementations:

```
module A1(A(..)) where
  data A = A { foo :: Int, bar :: Bool }

module A2(A(..)) where
  data A = A { foo :: Int, bar :: Bool }

module A(foo, bar) where
  import A1(foo)
  import A2(bar)
```

Here, module `A1` implements `signature A1`, module `A2` implements `signature A2`, and module `A` implements `signature A1` and `signature A2` individually and should certainly implement their merge. This is why we cannot simply merge type constructors based on the *OccName* of their top-level type; merging only occurs between in-scope identifiers.

Does merging a selector merge the type constructor?

```
signature A1(A(..)) where
  data A = A { foo :: Int, bar :: Bool }
```

```
signature A2(A(..)) where
  data A = A { foo :: Int, bar :: Bool }

signature A2(foo) where
  import A1(foo)
```

Does the last signature, which is written in the style of a sharing constraint on `foo`, also cause `bar` and the type and constructor `A` to be unified? Because a merge of a child name results in a substitution on the parent name, the answer is yes.

Incomplete data declarations

```
signature A1(A(foo)) where
  data A = A { foo :: Int }

signature A2(A(bar)) where
  data A = A { bar :: Bool }
```

Should `A1` and `A2` merge? If yes, this would imply that data definitions in signatures could only be *partial* specifications of their true data types. This seems complicated, which suggests this should not be supported; however, in fact, this sort of definition, while disallowed during type checking, should be *allowed* during shaping. The reason that the shape we ascribe to the signatures `A1` and `A2` are equivalent to the shapes for these which should merge:

```
signature A1(A(foo)) where
  data A = A { foo :: Int, bar :: Bool }

signature A2(A(bar)) where
  data A = A { foo :: Int, bar :: Bool }
```

4.3 Subtyping record selectors as functions

```
signature H(A, foo) where
  data A
  foo :: A -> Int

module M(A, foo) where
  data A = A { foo :: Int, bar :: Bool }
```

Does `M` successfully fill `H`? If so, it means that anywhere a signature requests a function `foo`, we can instead validly provide a record selector. This capability seems quite attractive, although in practice record selectors rarely seem to be abstracted this way: one reason is that `M.foo` still *is* a record selector, and can be used to modify a record. (Many library authors find this surprising!)

Nor does this seem to be an insurmountable instance of the avoidance problem: as a workaround, `H` can equivalently be written as:

```
signature H(foo) where
  data A = A { foo :: Int, bar :: Bool }
```

However, you might not like this, as the otherwise irrelevant `bar` must be mentioned in the definition.

In any case, actually implementing this ‘subtyping’ is quite complicated, because we can no longer assume that every child name is associated with a parent name. The technical difficulty is that we now need to unify a plain identifier *AvailInfo* (from the signature) with a type constructor *AvailInfo* (from a module.) It is not clear what this should mean. Consider this situation:

```

package p where
  signature H(A, foo, bar) where
    data A
    foo :: A -> Int
    bar :: A -> Bool
  module X(A, foo) where
    import H
package q where
  include p
  signature H(bar) where
    data A = A { foo :: Int, bar :: Bool }
  module Y where
    import X(A(..)) -- ???

```

Should the wildcard import on X be allowed? This question is equivalent to whether or not shaping discovers whether or not a function is a record selector and propagates this information elsewhere. If the wildcard is not allowed, here is another situation:

```

package p where
  -- define without record selectors
  signature X1(A, foo) where
    data A
    foo :: A -> Int
  module M1(A, foo) where
    import X1

package q where
  -- define with record selectors (X1s unify)
  signature X1(A(..)) where
    data A = A { foo :: Int, bar :: Bool }
  signature X2(A(..)) where
    data A = A { foo :: Int, bar :: Bool }

  -- export some record selectors
  signature Y1(bar) where
    import X1
  signature Y2(bar) where
    import X2

package r where
  include p
  include q

  -- sharing constraint
  signature Y2(bar) where
    import Y1(bar)

  -- the payload
  module Test where
    import M1(foo)
    import X2(foo)
    ... foo ... -- conflict?

```

Without the sharing constraint, the `foos` from `M1` and `X2` should conflict. With it, however, we should conclude that the `foos` are the same, even though the `foo` from `M1` is *not* considered a child of `A`, and even though in the sharing constraint we *only* unified `bar` (and its parent `A`). To know that `foo` from `M1` should also be unified, we have to know a bit more about `A` when the sharing constraint performs unification; however, the `AvailInfo` will only tell us about what is in-scope, which is *not* enough information.

5 Type checking

```

PkgType ::= ModIface0; ...; ModIfacen

Module interface
ModIface ::= module Module (mi_exports) where
                mi_decls
                mi_insts
                dep_orphs

mi_exports ::= AvailInfo0, ..., AvailInfon           Export list
mi_decls   ::= IfaceDecl0; ...; IfaceDecln           Defined declarations
mi_insts   ::= IfaceClsInst0; ...; IfaceClsInstn       Defined instances
dep_orphs  ::= Module0; ...; Modulen                 Transitive orphan dependencies

Interface declarations
IfaceDecl ::= OccName :: IfaceId
                | data OccName = IfaceData
                | ...
IfaceClsInst A type-class instance
IfaceId       Interface of top-level binder
IfaceData    Interface of type constructor

```

Figure 4: Module interfaces in GHC

In general terms, type checking an indefinite package (a package with holes) involves calculating, for every module, a `ModIface` representing the type/interface of the module in question (which is serialized to disk). The general form of these interface files are described in Figure 4; notably, the interfaces `IfaceId`, `IfaceData`, etc. contain `Name` references, which must be resolved by looking up a `ModIface` corresponding to the `Module` associated with the `Name`. (We will say more about this lookup process shortly.) For example, given:

```

package p where
  signature H where
    data T
  module A(S, T) where
    import H
    data S = S T

```

the `PkgType` is:

```

module HOLE:H (HOLE:H.T) where
  data T -- abstract type constructor
module THIS:A (THIS:A.S, HOLE:H.T) where
  data S = S HOLE:H.T
-- where THIS = p(H -> HOLE:H)

```

However, while it is true that the *ModIface* is the final result of type checking, we actually are conflating two distinct concepts: the user-visible notion of a *ModuleName*, which, when imported, brings some *Names* into scope (or could trigger a deprecation warning, or pull in some orphan instances...), versus the actual declarations, which, while recorded in the *ModIface*, have an independent existence: even if a declaration is not visible for an import, we may internally refer to its *Name*, and need to look it up to find out type information. (A simple case when this can occur is if a module exports a function with type $T \rightarrow T$, but doesn't export T).

$$\begin{aligned} \text{ModDetails} & ::= \langle \text{md_types}; \text{md_insts} \rangle \\ \text{md_types} & ::= \text{TyThing}_0, \dots, \text{TyThing}_n \\ \text{md_insts} & ::= \text{ClsInst}_0, \dots, \text{ClsInst}_n \end{aligned}$$

Type-checked declarations

<i>TyThing</i>	Type-checked thing with a <i>Name</i>
<i>ClsInst</i>	Type-checked type class instance

Figure 5: Semantic objects in GHC

Thus, a *ModIface* can be type-checked into a *ModDetails*, described in Figure 5. Notice that a *ModDetails* is just a bag of type-checkable entities which GHC knows about. We define the *external package state (EPT)* to simply be the union of the *ModDetails* of all external modules.

Type checking is a delicate balancing act between module interfaces and our semantic objects. A *ModIface* may get type-checked multiple times with different hole instantiations to provide multiple *ModDetails*. Furthermore complicating matters is that GHC does this resolution *lazily*: a *ModIface* is only converted to a *ModDetails* when we are looking up the type of a *Name* that is described by the interface; thus, unlike usual theoretical treatments of type checking, we can't eagerly go ahead and perform substitutions on *ModIfaces* when they get included.

In a separate compiler like GHC, there are two primary functions we must provide:

ModuleName to ModIface Given a *ModuleName* which was explicitly imported by a user, we must produce a *ModIface* that, among other things, specifies what *Names* are brought into scope. This is used by the renamer to resolve plain references to identifiers to real *Names*. (By the way, if shaping produced renamed trees, it would not be necessary to do this step!)

Module to ModDetails/EPT Given a *Module* which may be a part of a *Name*, we must be able to type check it into a *ModDetails* (usually by reading and typechecking the *ModIface* associated with the *Module*, but this process is involved). This is used by the type checker to find out type information on things.

There are two points in the type checker where these capabilities are exercised:

Source-level imports When a user explicitly imports a module, the *ModuleName* is mapped to a *ModIface* to find out what exports are brought into scope (*mi_exports*) and what orphan instances must be loaded (*dep_orphs*). Additionally, the *Module* is loaded to the EPT to bring instances from the module into scope.

Internal name lookup During type checking, we may have a *Name* for which we need type information (*TyThing*). If it's not already in the EPT, we type check and load into the EPT the *ModDetails* of the *Module* in the *Name*, and then check the EPT again. (`importDecl`)

5.1 *ModName* to *ModIface*

In all cases, the *mi_exports* can be calculated directly from the shaping process, which specifies exactly for each *ModName* in scope what will be brought into scope.

Does hiding a signature hide its orphans. Suppose that we have extended Backpack to allow hiding signatures from import.

```
package p requires (H) where -- H is hidden from import
  module A where
    instance Eq (a -> b) where -- orphan
    signature H {-# DEPRECATED "Don't use me" #-} where
    import A

package q where
  include p
  signature H where
    data T
  module M where
    import H -- warn deprecated?
    instance Eq (a -> b) -- overlap?
```

It is probably the most consistent to not pull in orphan instances and not give the deprecated warning: this corresponds to merging visible *ModIfaces*, and ignoring invisible ones.

Modules Modules are straightforward, as for any *Module* there is only one possibly *ModIface* associated with it (the *ModIface* for when we type-checked the (unique) `module` declaration.)

Signatures For signatures, there may be multiple *ModIfaces* associated with a *ModName* in scope, e.g. in this situation:

```
package p where
  signature S where
    data A
package q where
  include p
  signature S where
    data B
  module M where
    import S
```

Each literal `signature` has a *ModIface* associated with it; and the import of `S` in `M`, we want to see the *merged ModIfaces*. We can determine the *mi_exports* from the shape, but we also need to pull in orphan instances for each signature, and produce a warning for each deprecated signature.

5.2 *Module to ModDetails*

Modules For modules, we have a *Module* of the form $p(m \rightarrow Module, \dots)$, and we also have a unique *ModIface*, where each hole instantiation is `HOLE:m`.

To generate the *ModDetails* associated with the specific instantiation, we have to type-check the *ModIface* with the following adjustments:

1. Perform a *Module* substitution according to the instantiation of the *ModIface*'s *Module*. (NB: we *do* substitute `HOLE:A.x` to `HOLE:B.x` if we instantiated `A -> HOLE:B`, *unlike* the disjoint substitutions applied by shaping.)
2. Perform a *Name* substitution as follows: for any name with a package key that is a `HOLE`, substitute with the recorded *Name* in the requirements of the shape. Otherwise, look up the (unique) *ModIface* for the *Module*, and substitute with the corresponding *Name* in the *mi_exports*.

Signatures For signatures, we have a *Module* of the form `HOLE:m`. Unlike modules, there are multiple *ModIfaces* associated with a hole. We distinguish each separate *ModIface* by considering the full *PkgKey* it was defined in, e.g. `p(A -> HOLE:C, B -> q():B)`; call this the hole’s *defining package key*; the set of *ModIfaces* for a hole and their defining package keys can easily be calculated during shaping.

To generate the *ModDetails* associated with a hole, we type-check each *ModIface*, with the following adjustments:

1. Perform a *Module* substitution according to the instantiation of the defining package key. (NB: This may rename the hole itself!)
2. Perform a *Name* substitution as follows, in the same manner as would be done in the case of modules.
3. When these *ModDetails* are merged into the EPT, some merging of duplicate types may occur; a type may be defined multiple times, in which case we check that each definition is compatible with the previous ones. A concrete type is always compatible with an abstract type.

Invariants When we perform *Name* substitutions, we must be sure that we can always find out the correct *Name* to substitute to. This isn’t obviously true, consider:

```

package p where
  signature S(foo) where
    data T
    foo :: T
  module M(bar) where
    import S
    bar = foo
package q where
  module A(T(..)) where
    data T = T
    foo = T
  module S(foo) where
    import A
  include p
  module A where
    import M
    ... bar ...

```

When we type check `p`, we get the *ModIfaces*:

```

module HOLE:S(HOLE:S.foo) where
  data T
  foo :: HOLE:S.T
module THIS:M(THIS:M.bar) where
  bar :: HOLE:S.T

```

Now, when we type check `A`, we pull on the *Name* `p(S -> q():S):M.bar`, which means we have to type check the *ModIface* for `p(S -> q():S):M`. The un-substituted type of `bar` has a reference to `HOLE:S.T`; this should be substituted to `q():S.T`. But how do we discover this? We know that `HOLE:S` was instantiated to `q():S`, so we might try and look for `q():S.T`. However, this *Name* does not exist because the `module S` reexports the selector from `A`! Nor can we consult the (unique) *ModIface* for the module, as it doesn’t reexport the relevant type.

The conclusion, then, is that a module written this way should be disallowed. Specifically, the correctness condition for a signature is this: *Any Name mentioned in the ModIface of a signature must either be from an external module, or be exported by the signature.*

Special case export rule for record selectors. Here is the analogous case for record selectors:

```
package p where
  signature S(foo) where
    data T = T { foo :: Int }
  module M(bar) where
    import S
    bar = foo
package q where
  module A(T(..)) where
    data T = T { foo :: Int }
  module S(foo) where
    import A
  include p
  module A where
    import M
    ... bar ...
```

We could reject this, but technically we can find the right substitution for T, because the export of `foo` is an *AvailTC* which does mention T.

6 Cabal

Design goals:

- Backpack files are user-written. (In an earlier design, we had the idea that Cabal would generate Backpack files; however, we've since made Backpack files more user-friendly and reasonable to write by hand since they are reasonably designed for user development.)
- Backpack files are optional. A package can add a Backpack file to replace some (but not all) of the fields in a Cabal description.
- Backpack files can be compiled without GHC, if it is self-contained with respect to all the indefinite packages it includes. To include an indefinite package which is not locally defined but installed to the package database, you must use Cabal.
- Backpack packages are *unversioned*; you never see a version number in a Backpack package.

6.1 Versioning

In this section, we discuss how version numbers from Cabal factor into Backpack. In particular, versioning impacts the specification of *PkgKeys*. See <https://ghc.haskell.org/trac/ghc/wiki/Commentary/Packages/Concepts> for more background, and <https://ghc.haskell.org/trac/ghc/ticket/10566> for implementation progress.

Design goals Here are some design goals for versioning:

1. GHC doesn't know anything about version numbers: this is Cabal specific information. There are a few cases in GHC today where this design goal is already in force: pre-7.10, linker symbols were prefixed using a package name and version, but GHC simply represented this internally as an opaque string. And in today's GHC, package qualified imports only allow qualification by package name, and not by version.

2. Cabal doesn't know anything about package keys: GHC is responsible for calculating the package key of a package. This is because GHC must be able to maintain a mapping between the unhashed and hashed versions of a key, and the hashing process must be deterministic. If Cabal needs to generate a new package key, it must do so through GHC. (This is NOT how this is happening in GHC 7.10.)
3. Our design should, in principle, support mutual recursion between packages, even if the implementation does not at the moment.
4. GHC should not lose functionality, i.e. it should still be possible to link together the same package with different versions; however, Cabal may arrange for this to not occur by default unless a user explicitly asks for it.
5. A Cabal source package identifier (e.g. `foo-0.1`), which is a unit of distribution, is a distinct concept from a Backpack package (which we have referred to previously in the document as a mere package name), because a single Cabal file may ship a Backpack file that defines multiple internal packages.

These goals imply a few things:

1. Backpack files should not contain any version numbers, and should be agnostic to versioning. Backpack files are parsed and interpreted by GHC, and version numbers are Cabal's provenance!
2. As a corollary, if you want to refer to a specific version of a package from a Backpack file, this has to be done by giving the alternate version a different package name, e.g. `network-old`. (It is tempting to want to simply say that this means we should allow version numbers into GHC, but consider more complicated situations where you want to refer to two instances of `foo`, but one compiled with `bar-0.1` and the other compiled with `bar-0.2`, then your description of which package to pick up becomes considerably more complicated than just a package name and version. Better to defer this decision to Cabal.)
3. Package keys must record versioning information, otherwise we can't link together two different versions of the same package. This is due to our backwards-compatibility requirement.

Package keys To allow linking together multiple versions of the same package, we must record versioning information into the *PkgKey*. To do this, we include in the *PkgKey* a *VersionHash*. Cabal is responsible for defining *VersionHash* and may do whatever it wants, but we give two possible definitions in Figure 6.

p	Package name	
<i>SrcPkgId</i>	Cabal source package ID, e.g. <code>foo-0.1</code>	
<i>VersionHash</i>	$::= SrcPkgId \{ p_0 \rightarrow VersionHash_0, \dots p_n \rightarrow VersionHash_n \}$	Full version hash
<i>VersionHash'</i>	$::= SrcPkgId \{ SrcPkgId_0, \dots, SrcPkgId_n \}$	Simplified version hash
<i>PkgKey</i>	$::= p-VersionHash(m \rightarrow Module, \dots)$	

Figure 6: Version hash

The difference between a full version hash and a simplified version hash is what linking restrictions they impose on programs: the full version hash supports linking arbitrary versions of packages with arbitrary other versions, whereas the simplified hash has a Cabal-style requirement that there be some globally consistent mapping from package name to version.

The full version hash has some subtleties:

- Each sub-*VersionHash* recorded in a *VersionHash* is identified by a package name, which may not necessarily equal the package name embedded in the *SrcPkgId* in the *VersionHash*. This permits us to calculate a *VersionHash* for a package like:

```

package p where
  include network (Network)
  include network-old (Network as Network.Old)
  ...

```

if we want `network` to refer to `network-2.0` and `network-old` to refer to `network-1.0`. Without identifying each subdependency by package name, we could not distinguish the recorded *VersionHash*s for `network-old` and `network`.

- If a package name is locally specified in a Backpack file, it does not occur in the *VersionHash*: *VersionHash* strictly operates over Cabal’s notion of package identity.
- You might wonder why we need a *VersionHash* as well as a *PkgKey*; why not just specify *PkgKey* as *SrcPkgId* $\{p \rightarrow PkgKey, \dots\} (m \rightarrow Module, \dots)$? However, there is “too much” information in the *PkgKey*, causing the scheme to not work with mutual recursion:

```

package p where
  module M
  include q

```

To specify the package key of `p`, we need the package key of `q`; to specify the package key of `q`, we need the module identifier of `M` which contains the package key of `p`: circularity! (The simplified version hash does not have this problem as it is not recursive.)

6.2 Distribution and installation

How are Backpack files installed so other people can use them?

Challenges

- Prior to Backpack, when a Cabal package (e.g. unit of distribution) was compiled and installed would result in a single entry in the installed package database. With Backpack, compiling a package could result in multiple entries in the installed package database: (1) for indefinite packages which were instantiated, and (2) when there are multiple packages in a Backpack file.
- Relatedly, when we include an indefinite package, we may need to rebuild it with our specific dependencies. This makes compiling a Backpack file much more similar to `cabal-install` than to `Cabal`; however, the dependency structure is something that only GHC can calculate.

Why distribute Backpack files? Backpack files offer a convenient mechanism of defining multiple packages with inline syntax for modules. Further syntax extensions could allow us to give people a MixML style of programming in Haskell.

A Backpack file is not a replacement for a Cabal file: `exposed-modules` and similar fields are not necessary but we still need a `build-depends` to provide version bounds (until Backpack can also be used to handle version dependency.) This makes it easy for `cabal-install` to do its job.

This means we distinguish a package name `p` which occurs in a Backpack file and a Cabal *SrcPkgId*: Cabal creates a mapping between these. So to refer to an old version of a package, you would refer to it with a different name `q`, and then tell Cabal about the version bound constraints you want.

Definite packages Suppose we have written a Backpack file that looks like:

```

package helper where
  include base

```

```

    module P
package mypackage where
    include containers
    include helper
    module Q

```

and have written a Cabal file for it intending to distribute it on Hackage under the name `mypackage-0.1`. In the end, we will end up with the following entries in our installed package database:

```

name: "mypackage"
id: mypackage-1.0-IPID
version: 1.0
key: XXX
# e.g. mypackage-AAA {}
version-hash: AAA
# e.g. mypackage-1.0 { base -> base-4.7 , containers -> containers-0.5 }
depends: mypackage$helper-1.0-IPID, base-4.7-IPID
---
name: "mypackage$helper"
version: 1.0
id: mypackage$helper-1.0-IPID
key: YYY
# e.g. helper-AAA {}
version-hash: AAA
depends: containers-0.5-IPID

```

Things to note:

1. The package in the Backpack file with the same name as the Cabal package has special status: this is the package which is registered to the installed package database under the same name. All other packages are *qualified* under the Cabal package name, e.g. `mypackage$helper`.
2. The version hash, as described previously, is computed once for all packages in the Backpack file, and the `version` and `version-hash` are the same across all of them.
3. The key varies between the packages, since the `p` parameter is different in each one.
4. The installed package ID incorporates information about the package name.
5. Dependencies are only recorded directly **included** packages in a Backpack package. (GHC has to communicate to Cabal what the includes of every subpackage are.)

A more complex example with instantiated packages looks similar:

```

package helper where
    signature Data.Map
    module P
package mypackage where
    include containers (Data.Map)
    include helper
    module Q

```

however, now the instantiation is recorded in the database as well.

```

name: "mypackage"

```

```

id: mypackage-1.0-IPID
version: 1.0
key: XXX
# e.g. mypackage-AAA {}
version-hash: AAA
# e.g. mypackage-1.0 { containers -> containers-0.5 }
depends: mypackage$helper-1.0-IPID, containers-0.5-IPID
---
name: "mypackage$helper"
version: 1.0
id: mypackage$helper-1.0-IPID
key: YYY
# e.g. helper-AAA { Data.Map -> containers-KEY:Data.Map }
version-hash: AAA
depends: (none)
instantiated-with:
    Data.Map -> Data.Map@containers-0.5-IPID

```

More remarks:

1. Cabal's recorded `instantiated-with` records installed package IDs, so that the used implementation is uniquely determined.
2. Conversely, `depends` does NOT record non-textual dependencies such as instantiated holes. **is this necessary**
3. IPID includes information about how holes were instantiated.

GHC to Cabal When GHC compiles a Backpack file, it is the only entity which knows about the sub-packages of a package. In order to make sure they are all correctly installed, GHC has to communicate back some meta-data to Cabal: for each package,

- The (computed) package keys
- The dependencies
- The instantiation

I guess we have to define some format to do this. GHC can't directly write to the package database, because it doesn't know how to write in the Cabal-specific portion of the information.

This is clunky, is there a way to eliminate this? It's not possible for Cabal out of the box to handle this, since it assumes no module name conflicts but there definitely may be some in Backpack.

Indefinite package database The indefinite package database records indefinite packages (with holes) that have been typechecked. An indefinite package is associated with a (possibly unlimited) number of instantiated versions of the package, which have been fully instantiated and compiled.

An indefinite package is a new type of entry in the existing installed package database. **or maybe another entry in a different database** Here are the important things to keep track of for an indefinite package:

- Where do the (indefinite) interface files live? (NB: there are no libraries since we haven't compiled the package.)
- Where does the shape information live? (We could put it with the interface files, it's a pretty similar binary file.)

- Where does the source live, so we can recompile it when we instantiate it. (If it's empty, we'll have to refetch it from Hackage or something).
- Where does the Cabal configuration (result of running `cabal configure`) live, so that we build it with the same dependencies, flags, etc.

Associated with an indefinite package is some number of instantiated versions of this package. These are identified by package key (the installed package ID is the same) and are morally “sub”-packages of the indefinite package, although they get their own entries. **Alternate plan: put them together. Distinction between Cabal package and Backpack package.**

What makes installed indefinite packages difficult is that GHC may need to recompile them on the fly depending on an include.

The plan To be worked out