

Towards a more perfect union type

Functional pearl

Anonymous Author(s)

Abstract

We present a principled theoretical framework for dealing with union types, system, and show its work in practice on JSON data structures.

The framework poses union type inference as a problem of learning from multiple examples. Mathematical framework is quite generic, and easily extensible.

1 Introduction

Typing dynamic languages has been long considered a challenge [3]. The importance of the task grown with ubiquity cloud application programming interfaces (APIs) utilizing JavaScript object notation (JSON), where one needs to infer the structure having only a limited number of sample documents available.

Previous research have suggested it is possible to infer adequate type mappings from sample data [2, 10, 16].

In the present study, we expand on these results. We propose a framework for type systems in programming languages as learning algorithms, formulate it mathematically, and evaluate its performance on JSON API examples.

The proposed framework is grounded on mathematical theory, and complete typing relation. It is intended to add new features easily.

1.1 Related work

1.1.1 Union type providers

The earliest practical effort to apply union types to JSON inference to generate Haskell types[10]. It uses union type theory, but it also lacks an extensible theoretical framework.

F# type providers for JSON facilitate deriving a schema automatically; however, a type system is *ad-hoc*[16].

The other attempt to automatically infer schemas has been introduced in the PADS project [14]. Nevertheless, it has not specified a generalized type-system design methodology.

An approach presented with a program called [2] has been developed to derive types based on Markov chains. This approach requires considerable engineering time due to the implementation unit tests in a case-by-case mode, instead of formulating laws applying to all types. Moreover, this approach lacks sound underlying theory.

Therefore, we summarize that there are several previously introduced approaches that provide partially satisfactory results. In present study, we aim to expand these proposals to

enable systematic addition of features, and automatic validation of types.

1.1.2 Frameworks for describing type systems

Type systems are commonly expressed as partial relation of *typing*. Their properties, such as subject reduction are also expressed relatively to the relation (also partial) of *reduction* within a term rewriting system.

General formulations have been introduced for the Damas-Milner type systems parameterized by constraints [19].

We are not aware of any attempts to formulate general laws that would apply to all existing union type systems. Moreover, to the best of our knowledge no previous formulation exists that consider complete relations or functions in order to provide consistent mathematical descriptions where terms stray beyond their desired types¹.

It is also worth noting that traditional Damas-Milner type disciplines embrace the laws of soundness, and subject-reduction. However these laws often prove too strict during type system extension, and are abandoned in practice of larger systems[17].

2 Motivation

2.1 Motivating examples

Here, we consider several examples paraphrased from JSON API descriptions. These describe types underlying the motivation for the present study:

1. Subsets of data within a single constructor:

- *API argument is an email* – it is subset of valid String values, that can be validated on the client side.

```
{"example": [  
  "amy@example.com"  
  "robert@example.com"  
]}
```

```
example1a_values :: [Value]  
example1a_values = String <$> [  
  "amy@example.com"  
  , "edward@example.com"  
]  
example1a_repr :: HType  
example1a_repr = HRef "Email"
```

Haskell'20, August, 2020, Saint Louis, USA
2020.

¹Or at least beyond bottom expanding to *infamous undefined behaviour*[5].

- *The page size determines the number of results to return (min: 10, max:10,000)* - it is also a subset of integer values (Int) between 10, and 10,000

- *The date field contains ISO8601 date* – a record field is represented as a String that contains a calendar date in the format "2019-03-03"

2. Optional fields:

- *The page size is equal to 100 by default* - it means we have a record {"page_size": 50} or an empty record that should be interpreted as default value {}

3. Variant fields:

- *Answer to a query is either a number of registered objects, or String "unavailable"* - this is integer value (Int) or a String

4. Variant records:

- *Answer contains either a text message with an user id, or an error.* – That is can be represented as one of following options:

```
{ "message" : "Where can I submit my proposal?", "uid" : 1014 }
{ "message" : "Submit it to HotCRP", "uid" : 317 }
{ "error" : "Authorization failed", "code" : 401 }
{ "error" : "User not found", "code" : 404 }
```

5. Arrays corresponding to records²:

```
[
  [1, "Nick", null],
  [2, "George", "2019-04-11"],
  [3, "Olivia", "1984-05-03"]
]
```

6. Maps of identical objects³:

```
{
  "6408f5": {
    "size": 969709,
    "height": 510599,
    "difficulty": 866429.732,
    "previous": "54fced"
  },
  "54fced": {
    "size": 991394,
    "height": 510598,
    "difficulty": 866429.823,
    "previous": "6c9589"
  },
  "6c9589": {
    "size": 990527,
    "height": 510597,
    "difficulty": 866429.931,
    "previous": "51a0cb"
  }
}
```

²Which is considered a bad practice; however, it is part of real-life APIs. We may need to make it optional using the --array-records option.

³The example is taken from [2].

```
}
}
```

It should be noted that the last example presented above requires Haskell representation inference to be non-monotonic, as a dictionary with a single key would have an incompatible type:

```
data Example = Example { f_6408f5 :: O_6408f5 }
data O_6408f5 = O_6408f5 {
  size :: Int
, height :: Int
, difficulty :: Double
, previous :: String
}
```

It also suggests that a user might decide to explicitly add an evidence for one of alternative representations in the case when samples are insufficient. (like in case of a single element dictionary.)

2.2 Goal of inference

Given an undocumented (or incorrectly labelled) JSON API, we may need to read the input of Haskell encoding and avoid checking for the presence of *unexpected* format deviations. At the same time, we may decide to accept all known valid inputs outright so that we can use types⁴ to ensure that the input is processed exhaustively.

Accordingly, we can assume that the smallest non-singleton set is a better approximation type than a singleton set. We call it *minimal containing set principle*.

Second we can prefer types that allow for fewer number of *degrees of freedom* compared with the others, while conforming to a commonly occurring structure. We denote it as an *information content principle*.

Given these principles, and examples of frequently occurring patterns, we can infer a reasonable *world of types* that can be used as approximations, instead of establishing this procedure in an ad-hoc manner. In this way, we can implement *type system engineering*, that allows deriving type system design directly from the information about data structures and the likelihood of their occurrence.

3 Problem definition

3.1 Preliminaries

3.1.1 JSON values

As we focus on JSON, we utilize Haskell encoding of the JSON term for convenient reading⁵; specified as follows:

```
data Value =
  Object (Map String Value)
| Array [Value]
| String Text
| Number Scientific
```

⁴Compiler feature of checking for unmatched cases.

⁵As used by Aeson[1] package.

```

221 | Bool    Bool
222 | Null
223
224 To incorporate both integers and exact decimal fractions6
225 in the considered number representation, we employ deci-
226 mal floating point[18]:
227 data Scientific =
228   Scientific { coefficient :: Integer
229             , base10Exponent :: Int }

```

3.2 Defining type inference

3.2.1 Information in the type descriptions

If an inference fails, it is always possible to correct it by introducing an additional observation (example). To denote unification operation, or **information fusion** between two type descriptions, we use a Semigroup interface operation $\langle \rangle$ to merge types inferred from different observations.

```

238 class Semigroup ty where
239   (<>) :: ty -> ty -> ty

```

We use neutral element of the Monoid to indicate a type corresponding to no observations:

```

243 class Semigroup ty
244   => Monoid ty
245   where
246     mempty :: ty

```

In other words, we can say that mempty corresponds to situation when **no information was accepted** about a possible value (no term seen, not even a null). For example, an empty array [] can be referred to as an array type with mempty as an element type.

We describe the laws as QuickCheck [[@quickcheck](#)] properties so that unit testing can be implemented to detect obvious violations.

Neutral element of the Typelike monoid, mempty stands for **no information accepted** about possible value (no term seen, not even a null). For example an empty array [] could be typed as a array type with mempty as element type.

3.2.2 Beyond set

In the domain of permissive union types, a beyond set represents the case of **everything permitted** or a fully dynamic value, when we gather the information that permits every possible value inside a type. At the first reading, it may be deemed that a beyond set should comprise of only one single element – the top one.

However, since we defined **unification** operator $\langle \rangle$ as **information fusion**, we may encounter difficulties in assuring that no information has been lost during the unification⁷.

⁶JavaScript and JSON use a binary floating point instead; however we follow the representation selected by aeson library that parses JSON.

⁷Examples will be provided later.

Moreover, strict type systems usually specify more than one error value, as it should contain information about error messages and to keep track from where an error has been originated⁸.

This observation lets us consider type inference as a **learning problem**, and allows finding the common ground between the dynamic and the static typing discipline.

The languages relying on the static type discipline usually consider beyond as a set of error messages, as a value should correspond to a statically assigned and a **narrow** type. In this setting mempty as a fully polymorphic type forall a. a.

Languages with dynamic type discipline will treat beyond as untyped, dynamic value, and mempty again is a fully unknown, polymorphic value (like a type of an element of an empty array)⁹.

```

291 class (Monoid t, Eq t, Show t)
292   => Typelike t where
293   beyond :: t -> Bool

```

In addition, the standard laws for a **commutative** Monoid, we state the new law for the beyond set: The beyond set is always **closed to information addition** by $\langle \rangle a$ or $a \langle \rangle$ for any value of a. In other words the beyond set is an attractor of $\langle \rangle$ on both sides.¹⁰

Concerning union types, the key property of the beyond set, is that it is closed to information acquisition:

```

302 beyond_is_closed ty1 ty2 = do
303   beyond (ty1 :: ty) ==> beyond (ty1 <> ty2)

```

(We describe laws as QuickCheck properties~[4] so that unit testing can detect obvious violations.)

In this way, we can specify other elements of beyond set instead of a single top. When typing strict language, like Haskell, we seek to enable each element of the beyond set to contain at least one error message.¹¹

It should be noted that here, we abolish the semilattice requirement that has been conventionally assumed for type constraints [20], as this requirement is valid only for strict type constraint inference, not for a more general type inference considered as a learning problem. As we observe in the example [lst. 2.1](#), we need to perform non-monotonic inference when dealing with alternative representations.

When a specific instance of Typelike is also a semilattice (an idempotent semigroup), we will explicitly indicate if that is the case.

It is convenient validation when testing a recursive structure of the type.

Note that we abolish semilattice requirement that was traditionally assumed for type constraints here[21].

⁸7

⁹May sound similar until we consider adding more information to the type.

¹⁰So both in forall a. $\langle \rangle a$ and $\forall a. (a \langle \rangle)$ the result is kept in the beyond set.

¹¹It should be noted that many but not all type constraints are semilattice. Please refer to the counting example below.

That is because this requirement is valid only for strict type constraint inference, not for a more general type inference as a learning problem. As we saw in the example 1st. 2.1, we need non-monotonic inference when dealing with alternative representations.

It should be noted that this approach significantly generalized the assumptions compared with a full lattice subtyping [20][21].

3.2.3 Typing relation and its laws

The minimal definition of typing inference relation and type checking relation can be formulated as follows:

```
class Typelike ty
  => ty `Types` val where
  infer :: val -> ty
  check :: ty -> val -> Bool
```

Specifying the laws of typing is important, since we may need to consider separately the validity of a domain of types/type constraints, and that of the sound typing of the terms by these valid types.

First, we note that to describe *no information*, mempty cannot correctly type any term:

```
mempty_contains_no_terms term =
  check (mempty :: ty) term
  == False
```

Second important rule of typing is that all terms are typed successfully by any value in the beyond set.

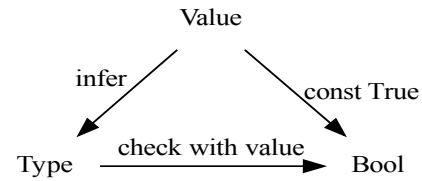
However, randomly drawing types for particular instances we might almost never get a type from the beyond set. In this case, we can use special generator called arbitraryBeyond that generates only the elements of the beyond set:

```
{-
beyond_contains_all_terms2 :: forall ty term.
  (Typelike ty
   ,Types ty term)
  => term -> _
beyond_contains_all_terms2 term =
  forAll arbitraryBeyond $ (`check` term)
-}
```

We state the most intuitive rule for typing: a type inferred from a term, must always be valid for that particular term.

```
inferred_type_contains_its_term ::
  forall ty term.
  ty `Types` term
  =>
  -> Bool
inferred_type_contains_its_term term =
  check ((infer:: term -> ty) term) (term :: term)
```

The law asserts that the following diagram commutes:



The last law states that the terms are correctly typechecked after adding more information into a single type. (For inference relation, it would be described as *principal type property*.)

```
fusion_keeps_terms v ty1 ty2 = do
  check ty1 v || check ty2 v ==>
  check (ty1 <> ty2) v
```

The minimal Typelike instance is the one that contains only mempty corresponding to the case of *no sample data received*, and a single beyond element for *all values permitted*. We will define it below as PresenceConstraint in sec. 3.5.

It should be noted that these laws are still compatible with the strict, static type discipline: namely the beyond set corresponds to a set of constraints with at least one type error, and a task of a compiler to prevent any program with the terms that type only to the beyond as a least upper bound.

3.3 Type engineering principles

Considering that we aim to infer a type from a finite number of samples, we are presented with a *learning problem*, so we need to use *prior* knowledge about the domain for inferring types.

Observing that $a : \text{false}$ we can expect that in particular cases, we may obtain that $a : \text{true}$. After noting that $b : 123$, we expect that $b : 100$ would also be acceptable. It means that we need to consider a typing system to *learn a reasonable general class from few instances*. This motivates formulating type system as an inference problem.

As the purpose is to deliver the most descriptive¹² types, we assume that we need to obtain a wider view rather than focusing on a *free type* and applying it to a larger sets whenever it is deemed justified.

The other principle corresponds to **correct operation**. It implies that having operations regarded on types, we can find a minimal set of types that assure correct operation on the case of unexpected errors.

Indeed we want to apply this theory to infer a type definition from a finite set of examples. We also seek to generalize it to infinite types.

For this purpose, we set the following rules of type design:

¹²The shortest one according to the information complexity principle.

- type should have a finite description
- inference must be a contravariant functor with regards to constructors. For example, if $\{ "a": X, "b": Y \}$ that is typed by $T \times y$, then $X :: x$ and $Y :: y$ must correspond to a valid typing.

3.3.1 Flat type constraints

Let us first consider typing of flat types: String and Number.

Constraints on number type First we infer the type description for integer values¹³:

```
data IntConstraint = IntRange Int Int
                  | IntNever
                  | IntAny
deriving (Show, Eq, Generic)
```

```
instance Semigroup IntConstraint where
  IntAny    <> _      = IntAny
  _         <> IntAny  = IntAny
  IntNever  <> a      = a
  a         <> IntNever = a
  IntRange  a      b <>
    IntRange c      d =
    IntRange (min a c) (max b d)
```

```
instance Typelike IntConstraint where
  beyond = (==IntAny)
```

```
instance Monoid IntConstraint where
  mempty = IntNever
```

```
instance IntConstraint `Types` Int where
  infer i = IntRange i i
  check IntNever _ = False
  check IntAny _ = True
  check (IntRange a b) i = a <= i && i <= b
```

JavaScript provides one number type that contains both Float and Int, so that the JSON values inherit this type:

```
data NumberConstraint =
  NCInt
  | NCNever
  | NCFloat
deriving(Eq,Show,Generic)
```

```
instance Semigroup NumberConstraint where
  NCFloat <> _      = NCFloat
  _       <> NCFloat = NCFloat
  NCNever <> a      = a
  a       <> NCNever = a
  NCInt   <> NCInt  = NCInt
```

¹³The implementation will make it optional with `--infer-int-ranges`.

```
instance Typelike NumberConstraint where
  beyond = (==NCFloat)
```

```
instance NumberConstraint `Types` Scientific where
  infer sci
    | base10Exponent sci >= 0 = NCInt
    | _                       = NCFloat
  check NCFloat _ = True
  check NCInt sci = base10Exponent sci >= 0
  check NCNever _ = False
```

```
instance Monoid NumberConstraint where
  mempty = NCNever
```

Constraints on string type

```
data StringConstraint =
  SCDate
  | SCEmail
  | SCEnum (Set Text)
  | SCNever
  | SCAny
deriving(Eq, Show, Generic)
```

```
instance StringConstraint `Types` Text where
  infer (isValidDate -> True) = SCDate
  infer (isValidEmail -> True) = SCEmail
  infer _ = SCAny
  infer value = SCEnum $ Set.singleton value
```

```
check SCDate s = isValidDate s
check SCEmail s = isValidEmail s
check (SCEnum vs) s = s `Set.member` vs
check SCNever _ = False
check SCAny _ = True
```

Then, whenever unifying the String constraint, the following code can be executed:

```
instance Semigroup StringConstraint where
  SCNever <> a      = a
  a       <> SCNever = a
  SCAny   <> _      = SCAny
  _       <> SCAny   = SCAny
  SCDate  <> SCDate  = SCDate
  SCEmail <> SCEmail = SCEmail
  (SCEnum a) <> (SCEnum b) |
    length (a <> b) < 10 = SCEnum (a <> b)
  _ <> _ = SCAny
```

```
instance Monoid StringConstraint where
  mappend = (<>)
  mempty = SCNever
```

```
instance Typelike StringConstraint where
  beyond = (==SCAny)
```

3.4 Free union type

Before we endeavour on finding type constraints for compound values (arrays and objects), it might be instructive to find a notion of *free type*, that is a type with no additional laws but the ones stated above.

Given a term with arbitrary constructors we can infer a *free type* for every term set T as follows: For any T value type $SetT$ satisfies our notion of *free type* specified as follows:

```
data FreeType a = FreeType { captured :: Set a }
  | Full
  deriving (Eq, Ord, Show, Generic)
```

```
instance (Ord a, Eq a) => Semigroup (FreeType a) where
  Full <> _ = Full
  _ <> Full = Full
  a <> b = FreeType
```

```
    $ (Set.union `on` captured) a b
instance (Ord a, Eq a) => Monoid (FreeType a) where
  mempty = FreeType Set.empty
```

```
instance (Ord a, Eq a, Show a)
  => Typelike (FreeType a) where
  beyond = (==Full)
```

```
instance (Ord a
  ,Eq a
  ,Show a)
  => FreeType a `Types` a where
  infer = FreeType . Set.singleton
  check Full _term = True
  check (FreeType s) term = term `Set.member` s
```

This definition is deemed sound, and may be applicable to a finite sets of terms or values. For a set of values: ["yes", "no", "error"], we may reasonably consider that type is an appropriate approximation of C-style enumeration, or Haskell-style ADT without constructor arguments.

However, the deficiency of this notion of *free type* is that it does not allow generalizing in infinite and recursive domains! It only allows utilizing objects from the sample.

3.5 Presence and absence constraint

We call this useful case a *presence or absence constraint*:

```
type role PresenceConstraint nominal
```

```
data PresenceConstraint a =
  Present
  | Absent
  deriving (Eq, Show, Typeable, Generic)
```

```
instance Semigroup (PresenceConstraint a) where
```

```
Absent <> a = a
a <> Absent = a
Present <> Present = Present
```

```
instance Monoid (PresenceConstraint a) where
  mempty = Absent
```

```
instance Typelike (PresenceConstraint a) where
  beyond = (==Present)
```

```
instance PresenceConstraint a `Types` a where
  infer _ = Present
  check Present _ = True
  check Absent _ = False
```

Although it does not seem useful in the context implying that we always have at least one input value, it is important as it can be used to specify an empty array (and therefore, an element type for which we observed no values).

After seeing true value we also expect false, so we can say that the basic constraint for a boolean value is its presence or absence.

```
type BoolConstraint = PresenceConstraint Bool
```

Note that booleans and null values are both denoted by this trivial PresenceConstraint constraint. The same is valid for null values, as there is only one null value.

```
type NullConstraint = PresenceConstraint ()
```

Note that we treat null as separate basic types, and postpone treatment of the union til later.

Variants Variants of two mutually exclusive types are also simple. They can be implement them with a type related to Either type that assumes these types are exclusive:

```
data a :| b = AltLeft a
  | AltRight b
  deriving (Show, Eq, Generic)
```

```
instance (FromJSON a
  ,FromJSON b)
  => FromJSON (a :| b) where
  parseJSON a = AltLeft <$> decodeEither
    <|> AltRight <$> decodeEither
```

In other words for Int :| String type, we first control whether the value is a String, and if this check fails, we attempt to parse it as String.

Variant records are slightly more complicated, as it may be unclear which typing is better to use:

```
{"message": "Where can I submit my proposal?",
  "uid" : 1014}
{"error" : "Authorization failed",
  "code" : 401}
```

```
data OurRecord =
  OurRecord { message :: Maybe String
```

```

661         , error  :: Maybe String
662         , code   :: Maybe Int
663         , uid    :: Maybe Int }
664
665     Or:
666     data OurRecord2 = Message { message :: String
667                               , uid     :: Int }
668                       | Error  { error  :: String
669                               , code   :: Int }

```

The best attempt here is to rely on the available examples being reasonably exhaustive. That is, we can estimate how many examples we have for each, and how many of them match. Then, we compare this number with type complexity (with options being more complex to process, because they need additional case expression.) In such cases, the latter definition has only one choice (optionality), but we only have two samples to begin with so we cannot be sure.

In the case of having more samples, the pattern emerges:

```

679 { "error" : "Authorization failed",
680   "code" : 401 }
681 { "message": "Where can I submit my proposal?",
682   "uid" : 1014 }
683 { "message": "Sent it to HotCRP",
684   "uid" : 93 }
685 { "message": "Thanks!",
686   "uid" : 1014 }
687 { "error" : "Missing user",
688   "code" : 404 }

```

Type cost function Since we are interested in types with less complexity and less optionality, we will define cost function as follows:

```

693 newtype TyCost = TyCost Int
694   deriving (Eq, Ord, Show, Enum, Num)
695
696 class Typelike ty
697   => TypeCost ty where
698   typeCost :: ty -> TyCost
699   typeCost a | a == mempty = 0
700               | otherwise   = 1
701
702 instance Semigroup TyCost where (< >) = (+)
703
704 instance Monoid TyCost where mempty = 0

```

When presented with several alternate representations from the same set of observations, we will use this function to select the least complex representation of the type. For flat constraints as above, we infer that they offer no optionality when no observations occurred (cost of 0), otherwise the cost is 1.

Considering that types beyond are to be avoided, we can assign conceptual *infinity* to these values. For the implementation purposes we will represent it by the value so high,

that is unlikely to ever occur in practical types, but still small enough that we can add it without checking for overflow.

```

718 inf :: TyCost
719 inf = 1000000000

```

Type cost should be non-negative, and non-decreasing when we add new observations to the type.

3.5.1 Object constraint

To avoid information loss, a constraint for JSON object type is introduced in such a way to **simultaneously gather information** about representing it either as a Map, or a record.

The typing of Map would be specified as follows:

```

729 data MappingConstraint =
730   MappingConstraint {
731     keyConstraint :: StringConstraint
732   , valueConstraint :: UnionType
733   }
734   | MappingNever
735   deriving (Eq, Show, Generic, Typeable)

```

```

737 instance Monoid MappingConstraint where
738   mempty = MappingNever

```

```

740 instance Typelike MappingConstraint where
741   beyond MappingNever = False
742   beyond MappingConstraint {..} =
743     beyond keyConstraint
744     && beyond valueConstraint

```

```

746 instance Semigroup MappingConstraint where
747   MappingNever <> a = a
748   a <> MappingNever = a
749   a <> b = MappingConstraint {
750     keyConstraint =
751       ((<>) `on` keyConstraint ) a b
752   , valueConstraint =
753       ((<>) `on` valueConstraint) a b
754   }

```

```

756 instance MappingConstraint `Types`
757   Object where
758   infer obj =
759     MappingConstraint
760       (foldMap infer $ Map.keys obj)
761       (foldMap infer      obj)
762   check MappingNever _ = False
763   check MappingConstraint {..} obj =
764     all (check keyConstraint)
765       (Map.keys obj)
766     && all (check valueConstraint)
767       (Foldable.toList obj)

```

Cost of mapping representation is a sum of cost of its fields:

```

771 instance TypeCost MappingConstraint where
772   typeCost MappingNever = 0
773   typeCost MappingConstraint {..} =
774     typeCost keyConstraint
775     + typeCost valueConstraint
776
777   Separately, we acquire the information about a possible
778   typing of a JSON object as a record of values:
779
780 data RecordConstraint =
781   RCTop
782   | RCBottom
783   | RecordConstraint {
784     fields :: HashMap Text UnionType
785   } deriving (Show,Eq,Generic, Typeable)
786
787 instance Typelike RecordConstraint where
788   beyond = (==RCTop)
789
790 instance Semigroup RecordConstraint where
791   RCBottom <> a = a
792   a <> RCBottom = a
793   RCTop <> _ = RCTop
794   _ <> RCTop = RCTop
795   a <> b = RecordConstraint $
796     Map.unionWith (<>) (fields a)
797     (fields b)
798
799 instance Monoid RecordConstraint where
800   mempty = RCBottom
801
802 instance RecordConstraint `Types` Object
803   where
804     infer = RecordConstraint
805       . Map.fromList
806       . fmap (second infer)
807       . Map.toList
808     check RCTop _ = True
809     check RCBottom _ = False
810     -- FIXME: treat extra keys!!!
811     check rc obj
812       | all (`elem` Map.keys (fields rc))
813         (Map.keys obj) =
814       and $ Map.elems $
815       Map.intersectionWith check
816         (fields rc)
817         obj
818     check _ _ = False
819
820 instance TypeCost RecordConstraint where
821   typeCost RCBottom = 0
822   typeCost RCTop = inf
823   typeCost RecordConstraint { fields } =
824     Foldable.foldMap typeCost fields
825

```

Observing that the two abstract domains considered above are independent, we can store the information about both options separately in a record¹⁴ as follows:

```

826 data ObjectConstraint = ObjectConstraint {
827   mappingCase :: MappingConstraint
828   , recordCase :: RecordConstraint
829   }
830   | ObjectNever
831   deriving (Eq,Show,Generic)
832

```

```

833 instance Semigroup ObjectConstraint where
834   ObjectNever <> a = a
835   a <> ObjectNever = a
836   a <> b =
837     ObjectConstraint {
838       mappingCase =
839         ((<>) `on` mappingCase) a b
840     , recordCase =
841         ((<>) `on` recordCase) a b
842     }
843

```

```

844 instance Monoid ObjectConstraint where
845   mempty = ObjectNever
846

```

```

847 instance Typelike ObjectConstraint where
848   beyond ObjectNever = False
849   beyond ObjectConstraint {..} =
850     beyond mappingCase
851     && beyond recordCase
852

```

```

853 instance ObjectConstraint `Types` Object where
854   infer v = ObjectConstraint (infer v)
855     (infer v)
856   check ObjectNever _ = False
857   check ObjectConstraint {..} v =
858     check mappingCase v
859     && check recordCase v
860

```

It should be noted that this representation is similar to *intersection type*: any value that satisfies `ObjectConstraint` must conform to both `mappingCase`, and `recordCase`.

It should be noted that this *intersection approach* to address alternative union type representations benefits from *principal type property*, meaning that a principal type is used to simply acquire the information corresponding to different representations and handle it separately.

Since we plan to choose only one representation for the object, we can say that minimum cost of this type is a minimum of component costs:

```

871 instance TypeCost ObjectConstraint where
872   typeCost ObjectNever = 0
873   typeCost ObjectConstraint {..} =
874

```

¹⁴Choice of representation will be explained later. Here we only consider acquiring the information about possible values.


```
typeCost mappingCase `min`
typeCost recordCase
```

3.5.2 Array constraint

Similarly to the object type, `ArrayConstraint` is used to simultaneously obtain information about all possible representations of an array, including the following:

- an array of the same elements;
- a row with the type depending on a column.

We need to acquire the information for both alternatives separately, and then, to measure a relative likelihood of either cases, before mapping the union type to Haskell declaration.

Here, we specify the records for two different possible representations:

```
data ArrayConstraint = ArrayConstraint {
  arrayCase :: UnionType
, rowCase :: RowConstraint
}
| ArrayNever
deriving (Show, Eq, Generic)
```

```
instance Monoid ArrayConstraint where
  mempty = ArrayNever
```

```
instance Typelike ArrayConstraint where
  beyond ArrayNever = False
  beyond ArrayConstraint {..} =
    beyond arrayCase
    && beyond rowCase
```

```
instance Semigroup ArrayConstraint where
  ArrayNever <> a = a
  a <> ArrayNever = a
  a1 <> a2 =
    ArrayConstraint {
      arrayCase = (((<>) `on` arrayCase) a1 a2
    , rowCase = (((<>) `on` rowCase ) a1 a2
    }
```

```
<<row-constraint>>
```

```
instance ArrayConstraint `Types` Array
  where
    infer vs =
      ArrayConstraint
        (mconcat (infer <$>
          Foldable.toList vs))
        (infer vs)
    check ArrayNever vs = False
    check ArrayConstraint {..} vs =
      and (check arrayCase <$>
```

```
Foldable.toList vs)
&& check rowCase vs
```

For the arrays, we plan to again choose only one of possible representations, so the cost of optionality is the lesser of the costs of the representation-specific constraints:

```
instance TypeCost ArrayConstraint where
  typeCost ArrayNever = 0
  typeCost ArrayConstraint {..} =
    typeCost arrayCase `min`
    typeCost rowCase
```

3.5.3 Row constraint

A row constraint is valid only if there is the same number of entries in all rows, which is represented by escaping the beyond set whenever there is an uneven number of columns.

```
data RowConstraint =
  RowTop
| RowNever
| Row [UnionType]
deriving (Eq, Show, Generic)
```

```
instance Typelike RowConstraint where
  beyond = (==RowTop)
```

```
instance Monoid RowConstraint where
  mempty = RowNever
```

```
instance RowConstraint `Types` Array where
  infer = Row
    . Foldable.toList
    . fmap infer
  check RowTop _ = True
  check RowNever _ = False
  check (Row rs) vs
    | length rs == length vs =
      and $
        zipWith check rs
          (Foldable.toList vs)
  check _ _ = False
```

```
instance Semigroup RowConstraint where
  RowTop <> _ = RowTop
  _ <> RowTop = RowTop
  RowNever <> a = a
  a <> RowNever = a
  Row bs <> Row cs
    | length bs /= length cs = RowTop
  Row bs <> Row cs =
    Row $ zipWith (<>) bs cs
```

In other words, `RowConstraint` is a *levitated semilattice* with a neutral element [12] over the content type that is a list of `UnionType` objects.

The cost of the row constraint is inferred in a similar manner as the cost of the record constraint:

```
instance TypeCost RowConstraint where
  typeCost RowNever = 0
  typeCost RowTop   = inf
  typeCost (Row cols) = foldMap typeCost cols
```

3.5.4 Combining the above into a union type

It should note that given the constraints for the different type constructors, the union type can be considered as mostly a generic Monoid instance[8]:

```
data UnionType =
  UnionType {
    unionNull :: NullConstraint
  , unionBool :: BoolConstraint
  , unionNum  :: NumberConstraint
  , unionStr  :: StringConstraint
  , unionArr  :: ArrayConstraint
  , unionObj  :: ObjectConstraint
  } deriving (Eq, Generic)
```

```
instance Semigroup UnionType where
```

```
  u1 <> u2 =
    UnionType {
      unionNull = ((<>) `on` unionNull) u1 u2
    , unionBool = ((<>) `on` unionBool) u1 u2
    , unionNum  = ((<>) `on` unionNum ) u1 u2
    , unionStr  = ((<>) `on` unionStr ) u1 u2
    , unionObj  = ((<>) `on` unionObj ) u1 u2
    , unionArr  = ((<>) `on` unionArr ) u1 u2
    }
```

The generic structure of union type can be explained by the fact that the information contained in each record field is *independent* from the information contained in other fields. It means that we perform unification independently over different dimensions.

```
instance Monoid UnionType where
```

```
  mempty = UnionType {
    unionNull = mempty
  , unionBool = mempty
  , unionNum  = mempty
  , unionStr  = mempty
  , unionObj  = mempty
  , unionArr  = mempty
  }
```

As we described previously, the beyond set may correspond to either **accepting any value** or to **accepting no more information**. Its definition should be no surprise:

```
instance Typelike UnionType where
```

```
  beyond UnionType {..} =
    beyond unionNull
  && beyond unionBool
```

```
&& beyond unionNum
&& beyond unionStr
&& beyond unionObj
&& beyond unionArr
```

Inference breaks down disjoint alternatives corresponding to different record fields, depending on the constructor of a given value.

It enables implementing a clear and efficient treatment of different alternatives separately¹⁵.

```
instance UnionType `Types` Value where
  infer (Bool b) = mempty { unionBool = infer b }
  infer Null     = mempty { unionNull = infer () }
  infer (Number n) = mempty { unionNum  = infer n }
  infer (String s) = mempty { unionStr  = infer s }
  infer (Object o) = mempty { unionObj  = infer o }
  infer (Array a) = mempty { unionArr  = infer a }
  check UnionType { unionBool } (Bool b) =
    check unionBool b
  check UnionType { unionNull } Null =
    check unionNull ()
  check UnionType { unionNum } (Number n) =
    check unionNum n
  check UnionType { unionStr } (String s) =
    check unionStr s
  check UnionType { unionObj } (Object o) =
    check unionObj o
  check UnionType { unionArr } (Array a) =
    check unionArr a
```

Since union type is all about optionality, we need to sum all options from different alternatives:

```
instance TypeCost UnionType where
  typeCost UnionType {..} = typeCost unionBool
    + typeCost unionNull + typeCost unionNum
    + typeCost unionStr + typeCost unionObj
    + typeCost unionArr
```

3.5.5 Overlapping alternatives

The essence of union type systems have long been dealing with the conflicting types provided in the input.

Motivated by the examples above, we also aim to address conflicting alternative assignments.

It is apparent that examples 4. to 6. hint at more than one assignment:

5. A set of lists of values that may correspond to Int, String, or null, or a table that has the same (and predefined) type for each values.
6. A record of fixed names or the mapping from hash to a single object type.

¹⁵The question may arise: what is the *union type* without *set union*? When the sets are disjoint, we just put the values in different bins to enable easier handling.

3.5.6 Counting observations

In this section, we discuss how to gather information about the number of samples supporting each alternative type constraint. To explain this, the other example can be considered:

```
{ "samples":
  [{ "error" : "Authorization failed",
    "code" : 401}
  , { "message": "Where can I submit my proposal?",
    "uid" : 1014}
  , { "message": "Sent it to HotCRP",
    "uid" : 93}
  , { "message": "Thanks!",
    "uid" : 1014}
  , { "error" : "Authorization failed",
    "code": 401}
  ] }
```

First, we need to identify it as a list of similar elements. Second, we note, that there are multiple instances of each record example. We consider that the best approach would be to use the multisets of inferred records instead of normal sets. To find the best representation, we can a type complex-ity, and attempt to minimize the term.

Next step is to detect the similarities between type descriptions introduced for different parts of the term:

```
{ "samples"      : [...],
  "last_message" : { "message": "Thanks!",
                    "uid" : 1014}
}
```

We can add the auxiliary information about a number of samples observed, and the constraint remains a Typelike object:

```
data Counted a =
  Counted { count    :: Int
          , constraint :: a
          } deriving (Eq, Show, Generic)

instance Semigroup a
  => Semigroup (Counted a) where
  a <> b = Counted (count a + count b)
          (constraint a <> constraint b)

instance Monoid a
  => Monoid (Counted a) where
  mempty = Counted 0 mempty

instance Typelike a
  => Typelike (Counted a) where
  beyond Counted {..} = beyond constraint

instance ty `Types` term
  => (Counted ty) `Types` term where
  infer term = Counted 1 $ infer term
```

```
check (Counted _ ty) term = check ty term
```

```
instance TypeCost ty
  => TypeCost (Counted ty) where
  typeCost (Counted _ ty) = typeCost ty
```

We can interconnect Counted as parametric functor to select constraints to track auxiliary information.

It should be noted that Counted constraint is the first example that does not correspond to a semilattice, that is $a <> a \neq a$.

This is because it is a Typelike object; however, it is not a type constraint in a conventional sense. Instead it counts the number of samples observed for the constraint inside so that we can decide on which alternative representation is best supported by evidence.

Therefore, at each step, we may need to maintain a **cardinality** of each possible value, and being provided with sufficient number of samples, we may attempt to detect.¹⁶

To preserve efficiency, we may need to merge whenever the number of alternatives in a multiset crosses the threshold.¹⁷ We can attempt to narrow strings only in the cases when cardinality crosses the threshold.¹⁸

4 Selecting representations

4.1 Specifying heuristics to achieve better types

The final touch would be to perform the post-processing of an assigned type before generating it to make it more resilient to common uncertainties.

It should be noted that these assumptions may bypass the defined least-upper-bound criterion specified in the initial part of the paper; however, they prove to work well in practice.

4.1.1 Promoting empty type

If we have no observations corresponding to an array type, it can be inconvenient to disallow an array to contain any values at all. Therefore, we introduce a non-monotonic step of converting the mempty into a final Typelike object aiming to introduce a representation allowing the occurrence of any Value in the input. That still preserves the validity of the typing.

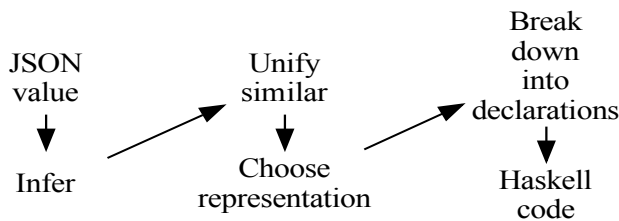
We note that the proposed program must not have any assumptions about these values; however, at the same time it should be able to print them for debugging purposes.

¹⁶If we detect a pattern too early, we risk to make the types too narrow to work with actual API answers.

¹⁷Option --max-alternative-constructors=N.

¹⁸Option --min-enumeration-cardinality.

4.2 Overall processing scheme



4.3 Simplification by identifying unification candidates

In most JSON documents, we observe that the same object can be described in different parts of sample data structures. Due to this reason, we compare the sets of labels assigned to all objects and propose to unify those that have more than 60% of identical labels.

For transparency, the identified candidates are logged for each user, and a user can also indicate them explicitly instead of relying on automation.

We conclude that this allows considerably decreasing the complexity of types and makes the output less redundant.

5 Future work

5.1 Scaling to type environments

In the present paper, we only discuss typing of tree-like values. However, it is natural to scale this approach to multiple types in APIs, in which different types are referred to by name and possibly contain each other.

To address these cases, we will show that the environment of Typelike objects is also Typelike, and that constraint unification can be extended in the same way.

5.2 Generic derivation of Typelike

It should be noted that Typelike instances for non-simple types usually follow one the two patterns:

1. for typing terms that have a finite sum of disjoint constructors, we bin this information by each constructor during the inference
2. for typing terms with multiple alternative representations, we infer all constraints separately for each representation by applying a different inference algorithm to the same term

In both cases, the derivation procedure of the Monoid, and Typelike instances is the same.

It allows using GHC Generics[13] to specify standard implementations for most of the boilerplate code.

It means that we only have to manually define the following:

- new constraint data types¹⁹,

¹⁹In many cases one can also rely on a generic constraint representation derived from Generic representation type Rep, that is when inference is mutually exclusive by term type constructors.

- inference from constructors (case 1), as well as providing the entirety of handling alternative constraints until we select representations.

5.3 Conclusion

In the present study, we aimed to derive the types that were valid with respect to the provided specification, thereby obtaining the information from the input in most comprehensive way.

We defined type inference as representation learning and type system engineering as a meta-learning problem in which the **priors corresponding to the data structure induced typing rules**.

We also formulated the **union type discipline** as manipulation of Typelike commutative monoids, that represented knowledge about the data structure.

In addition, we proposed a union type system engineering methodology that was logically justified by a theoretical criteria. We demonstrated that it was capable of consistently explaining the decisions made in practice.

We consider that this kind of *formally justified type system engineering* can become widely used in practice, replacing *ad-hoc* approaches in the future.

The proposed approach may be used to underlie the way towards formal construction and derivation of type systems based on the specification of value domains and design constraints.

Bibliography

- [1] Aeson: Fast JSON parsing and generation: 2011. <https://hackage.haskell.org/package/aeson>.
- [2] A first look at quicktype: 2017. <https://blog.quicktype.io/first-look/>.
- [3] Anderson, C. et al. 2005. Towards type inference for javascript. *ECOOP 2005 - object-oriented programming* (Berlin, Heidelberg, 2005), 428–452.
- [4] Claessen, K. and Hughes, J. 2000. QuickCheck: A lightweight tool for random testing of haskell programs. *SIGPLAN Not.* 35, 9 (Sep. 2000), 268–279. DOI:<https://doi.org/10.1145/357766.351266>.
- [5] C Standard undefined behaviour versus Wittgenstein: <https://www.yodaiken.com/2018/05/20/depressing-and-faintly-terrifying-days-for-the-c-standard/>.
- [6] C Undefined Behavior - Depressing and Terrifying (Updated): 2018. <https://www.yodaiken.com/2018/05/20/depressing-and-faintly-terrifying-days-for-the-c-standard/>.
- [7] EnTangleD: A bi-directional literate programming tool: 2019. <https://blog.esciencecenter.nl/entangled-1744448f4b9f>.
- [8] Generics example: Creating monoid instances: 2012. <https://www.yesodweb.com/blog/2012/10/generic-monoid>.

- [9] GHCID - a new ghci based ide (ish): 2014. <http://neilmitchell.blogspot.com/2014/09/ghcid-new-ghci-based-ide-ish.html>.
- [10] JSON autotype: Presentation for Haskell.SG: 2015. <https://engineers.sg/video/json-autotype-1-0-haskell-sg--429>.
- [11] Knuth, D.E. 1984. Literate programming. *Comput. J.* 27, 2 (May 1984), 97–111. DOI:<https://doi.org/10.1093/comjnl/27.2.97>.
- [12] Lattices: Fine-grained library for constructing and manipulating lattices: 2017. <http://hackage.haskell.org/package/lattices-2.0.2/docs/Algebra-Lattice-Levitated.html>.
- [13] Magalhães, J.P. et al. 2010. A generic deriving mechanism for haskell. *SIGPLAN Not.* 45, 11 (Sep. 2010), 37–48. DOI:<https://doi.org/10.1145/2088456.1863529>.
- [14] Mandelbaum, Y. et al. 2007. PADS/ML: A Functional Data Description Language. *Proceedings of the 34th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (New York, NY, USA, 2007), 77–83.
- [15] Pandoc: A universal document converter: 2000. <https://pandoc.org>.
- [16] Petricek, T. et al. 2016. Types from Data: Making Structured Data First-Class Citizens in F#. *SIGPLAN Not.* 51, 6 (Jun. 2016), 477–490. DOI:<https://doi.org/10.1145/2980983.2908115>.
- [17] Peyton Jones, S. 2019. Type inference as constraint solving: How ghc’s type inference engine actually works. Zurihac keynote talk.
- [18] scientific:Numbers represented using scientific notation: <https://hackage.haskell.org/package/scientific-0.3.6.2>.
- [19] Sulzmann, M. and Stuckey, P.j. 2008. Hm(x) Type Inference is Clp(x) Solving. *J. Funct. Program.* 18, 2 (Mar. 2008), 251–283. DOI:<https://doi.org/10.1017/S0956796807006569>.
- [20] Tiuryn, J. 1992. Subtype inequalities. [1992] *Proceedings of the Seventh Annual IEEE Symposium on Logic in Computer Science.* (1992), 308–315.
- [21] Tiuryn, J. 1997. Subtyping over a lattice (abstract). *Computational logic and proof theory* (Berlin, Heidelberg, 1997), 84–88.
- [22] Undefined behavior in 2017: 2017. <https://blog.regehr.org/archives/1520>.

Appendix: definition module headers

```
{-# language AllowAmbiguousTypes #-}
{-# language DeriveGeneric #-}
{-# language DuplicateRecordFields #-}
{-# language FlexibleInstances #-}
{-# language GeneralizedNewtypeDeriving #-}
{-# language MultiParamTypeClasses #-}
{-# language NamedFieldPuns #-}
{-# language PartialTypeSignatures #-}
{-# language ScopedTypeVariables #-}
{-# language TypeOperators #-}
```

```
{-# language RoleAnnotations #-}
{-# language ViewPatterns #-}
{-# language RecordWildCards #-}
{-# language OverloadedStrings #-}
{-# ghc_options -Wno-orphans #-}
module Unions where

import Control.Arrow(second)
import Data.Aeson
import Data.Maybe(isJust,catMaybes)
import qualified Data.Foldable as Foldable
import Data.Function(on)
import Data.Text(Text)
import qualified Data.Text as Text
import qualified Data.Text.Encoding as Text
import qualified Text.Email.Validate(isValid)
import qualified Data.Set as Set
import Data.Set(Set)
import Data.Scientific
import Data.String
--import Data.List(sortBy)
import qualified Data.HashMap.Strict as Map
import Data.HashMap.Strict(HashMap)
import GHC.Generics(Generic)
import Data.Hashable
import Data.Typeable
import Data.Time.Format(iso8601DateFormat,parseTimeM,default)
import Data.Time.Calendar(Day)
```

```
<<freetype>>
<<typelike>>
<<basic-constraints>>
<<array-constraint>>
<<object-constraint>>
<<presence-absence-constraints>>
<<union-type-instance>>
<<type>>
<<counted>>
<<typecost>>
<<representation>>

<<missing>>
```

Appendix: test suite

```
{-# language FlexibleInstances #-}
{-# language Rank2Types #-}
{-# language MultiParamTypeClasses #-}
{-# language NamedFieldPuns #-}
{-# language ScopedTypeVariables #-}
{-# language StandaloneDeriving #-}
{-# language TemplateHaskell #-}
{-# language TypeOperators #-}
```



```

1431 {-# language TypeApplications #-}
1432 {-# language TupleSections #-}
1433 {-# language UndecidableInstances #-}
1434 {-# language AllowAmbiguousTypes #-}
1435 {-# language OverloadedStrings #-}
1436 {-# language ViewPatterns #-}
1437 {-# ghc_options -Wno-orphans #-}
1438 module Main where
1439
1440 import qualified Data.Set as Set
1441 import qualified Data.Text as Text
1442 import qualified Data.ByteString.Char8 as BS
1443 import Control.Monad(when)
1444 import Data.FileEmbed
1445 import Data.Maybe
1446 import Data.Scientific
1447 import Data.Aeson
1448 import Data.Proxy
1449 import Data.Typeable
1450 import Test.Hspec
1451 import Test.Hspec.QuickCheck
1452 import Test.QuickCheck
1453 import Test.Validity.Shrinking.Property
1454 import Test.Validity.Utils(nameOf)
1455 import qualified GHC.Generics as Generic
1456 import Test.QuickCheck.Classes
1457 import System.Exit(exitFailure)
1458
1459 import Test.Arbitrary
1460 import Test.LessArbitrary as LessArbitrary
1461 import Unions
1462
1463 instance Arbitrary Value where
1464     arbitrary = fasterArbitrary
1465
1466 instance LessArbitrary Value where
1467     lessArbitrary = cheap $$$$ genericLessArbitrary
1468     where
1469         cheap = LessArbitrary.oneof [
1470             pure Null
1471             , Bool <$> lessArbitrary
1472             , Number <$> lessArbitrary
1473             ]
1474
1475 instance LessArbitrary a
1476     => LessArbitrary (Counted a) where
1477
1478 instance LessArbitrary a
1479     => Arbitrary (Counted a) where
1480     arbitrary = fasterArbitrary
1481
1482 instance Arbitrary Object where
1483     arbitrary = fasterArbitrary
1484
1485

```

```

1486 instance Arbitrary Array where
1487     arbitrary = fasterArbitrary
1488
1489 class Typelike ty
1490     => ArbitraryBeyond ty where
1491     arbitraryBeyond :: CostGen ty
1492
1493 instance ArbitraryBeyond (PresenceConstraint a) where
1494     arbitraryBeyond = pure Present
1495
1496 instance ArbitraryBeyond StringConstraint where
1497     arbitraryBeyond = pure SCAny
1498
1499 instance ArbitraryBeyond IntConstraint where
1500     arbitraryBeyond = pure IntAny
1501
1502 instance ArbitraryBeyond NumberConstraint where
1503     arbitraryBeyond = pure NCFloat
1504
1505 instance ArbitraryBeyond RowConstraint where
1506     arbitraryBeyond = pure RowTop
1507
1508 instance ArbitraryBeyond RecordConstraint where
1509     arbitraryBeyond = pure RCTop
1510
1511 instance ArbitraryBeyond MappingConstraint where
1512     arbitraryBeyond =
1513         MappingConstraint <$$$> arbitraryBeyond
1514         <*> arbitraryBeyond
1515
1516 instance (Ord a
1517         , Show a
1518         )
1519     => ArbitraryBeyond (FreeType a) where
1520     arbitraryBeyond = pure Full
1521
1522 instance ArbitraryBeyond ObjectConstraint where
1523     arbitraryBeyond = do
1524         ObjectConstraint <$$$> arbitraryBeyond
1525         <*> arbitraryBeyond
1526
1527 instance ArbitraryBeyond ArrayConstraint where
1528     arbitraryBeyond = do
1529         ArrayConstraint <$$$> arbitraryBeyond
1530         <*> arbitraryBeyond
1531
1532 instance ArbitraryBeyond UnionType where
1533     arbitraryBeyond =
1534         UnionType <$$$> arbitraryBeyond
1535         <*> arbitraryBeyond
1536         <*> arbitraryBeyond
1537         <*> arbitraryBeyond
1538         <*> arbitraryBeyond
1539
1540

```

```

1541
1542 arbitraryBeyondSpec :: forall      ty.
1543     (ArbitraryBeyond ty
1544      ,Typelike      ty)
1545     => Spec
1546 arbitraryBeyondSpec =
1547     prop "arbitrarybeyond returns terms beyond" $
1548     (beyond <$> (arbitraryBeyond :: CostGen ty))
1549
1550 instance LessArbitrary Text.Text where
1551     lessArbitrary = Text.pack <$> lessArbitrary
1552
1553 instance Arbitrary Text.Text where
1554     arbitrary = Text.pack <$> arbitrary
1555
1556 instance Arbitrary Scientific where
1557     arbitrary = scientific <$> arbitrary
1558               <*> arbitrary
1559
1560 instance (LessArbitrary      a
1561           ,Ord                a)
1562           => LessArbitrary (FreeType a) where
1563
1564 instance Arbitrary (FreeType Value) where
1565     arbitrary = fasterArbitrary
1566     {-shrink Full      = []
1567      shrink (FreeType elts) = map FreeType
1568                               $ shrink elts-}
1569
1570 instance (Ord                v
1571           ,Show              v)
1572           => TypeCost (FreeType v) where
1573     typeCost Full      = inf
1574     typeCost (FreeType s) = TyCost $ Set.size s
1575
1576 {-
1577 instance (Eq                a
1578           ,Ord              a
1579           ,GenUnchecked     a
1580           ,LessArbitrary    a
1581           ,LessArbitrary (FreeType a)
1582           ,Arbitrary        (FreeType a))
1583           => GenUnchecked (FreeType a) where
1584     genUnchecked = fasterArbitrary
1585     shrinkUnchecked Full      = []
1586     shrinkUnchecked FreeType { captured } =
1587         map (FreeType . Set.fromList)
1588             $ shrinkUnchecked
1589             $ Set.toList captured
1590
1591 instance Validity (FreeType a) where
1592     validate _ = validate True
1593 -}
1594
1595

```

```

instance LessArbitrary (PresenceConstraint a) where
    lessArbitrary = genericLessArbitraryMonoid
instance Arbitrary      (PresenceConstraint a) where
    arbitrary = fasterArbitrary
instance LessArbitrary IntConstraint where
    lessArbitrary = genericLessArbitraryMonoid
instance Arbitrary      IntConstraint where
    arbitrary = fasterArbitrary
instance LessArbitrary NumberConstraint where
    lessArbitrary = genericLessArbitrary
instance Arbitrary      NumberConstraint where
    arbitrary = fasterArbitrary
instance LessArbitrary StringConstraint where
    lessArbitrary = genericLessArbitraryMonoid
instance Arbitrary      StringConstraint where
    arbitrary = fasterArbitrary
instance LessArbitrary ObjectConstraint where
    lessArbitrary = genericLessArbitraryMonoid
instance Arbitrary      ObjectConstraint where
    arbitrary = fasterArbitrary
instance LessArbitrary RecordConstraint where
    lessArbitrary = genericLessArbitraryMonoid
instance Arbitrary      RecordConstraint where
    arbitrary = fasterArbitrary
instance LessArbitrary ArrayConstraint where
    lessArbitrary = genericLessArbitraryMonoid
instance Arbitrary      ArrayConstraint where
    arbitrary = fasterArbitrary
instance LessArbitrary RowConstraint where
    lessArbitrary = genericLessArbitraryMonoid
instance Arbitrary      RowConstraint where
    arbitrary = fasterArbitrary
instance LessArbitrary MappingConstraint where
    lessArbitrary = genericLessArbitraryMonoid
instance Arbitrary      MappingConstraint where
    arbitrary = fasterArbitrary
instance LessArbitrary UnionType where
    lessArbitrary = genericLessArbitraryMonoid
instance Arbitrary      UnionType where
    arbitrary = fasterArbitrary
{-

```

```

1651 instance GenUnchecked UnionType where
1652   genUnchecked = arbitrary
1653   shrinkUnchecked = shrink
1654 -}
1655 {-
1656 instance Validity UnionType where
1657   validate _ = validate True-}
1658
1659 shrinkSpec :: forall a.
1660   (Arbitrary a
1661    ,Typeable a
1662    ,Show a
1663    ,Eq a
1664    )
1665   => Spec
1666 shrinkSpec = prop ("shrink on " <> nameOf @a)
1667   $ doesNotShrinkToItself arbitrary (shrink :: a -> [a])
1668
1669 allSpec :: forall ty v.
1670   (Typeable ty
1671    ,Arbitrary ty
1672    ,Show ty
1673    ,Types ty v
1674    ,ArbitraryBeyond ty
1675    ,Arbitrary v
1676    ,Show v
1677    ) => Spec
1678 allSpec = describe (nameOf @ty) $ do
1679   arbitraryBeyondSpec @ty
1680   shrinkSpec @ty
1681
1682 <<typelike-spec>>
1683 <<types-spec>>
1684 <<typecost-laws>>
1685
1686 main :: IO ()
1687 main = do
1688   exitFailure
1689   putStrLn "NumberConstraint"
1690   {-sample $ arbitrary @Value
1691   sample $ arbitrary @NullConstraint
1692   sample $ arbitrary @NumberConstraint
1693   sample $ arbitrary @RowConstraint
1694   sample $ arbitrary @RecordConstraint
1695   sample $ arbitrary @ArrayConstraint
1696   sample $ arbitrary @MappingConstraint
1697   sample $ arbitrary @ObjectConstraint-}
1698
1699 lawsCheckMany
1700   [typesSpec (Proxy :: Proxy (FreeType Value) )
1701     (Proxy :: Proxy Value ) True
1702   ,typesSpec (Proxy :: Proxy NumberConstraint )
1703     (Proxy :: Proxy Scientific) True
1704   ]
1705 
```

```

1706 ,typesSpec (Proxy :: Proxy StringConstraint )
1707   (Proxy :: Proxy Text.Text ) True
1708 ,typesSpec (Proxy :: Proxy BoolConstraint )
1709   (Proxy :: Proxy Bool ) True
1710 ,typesSpec (Proxy :: Proxy NullConstraint )
1711   (Proxy :: Proxy () ) True
1712 ,typesSpec (Proxy :: Proxy RowConstraint )
1713   (Proxy :: Proxy Array ) True
1714 ,typesSpec (Proxy :: Proxy ArrayConstraint )
1715   (Proxy :: Proxy Array ) True
1716 ,typesSpec (Proxy :: Proxy MappingConstraint)
1717   (Proxy :: Proxy Object ) True
1718 ,typesSpec (Proxy :: Proxy RecordConstraint )
1719   (Proxy :: Proxy Object ) True
1720 ,typesSpec (Proxy :: Proxy ObjectConstraint )
1721   (Proxy :: Proxy Object ) True
1722 ,typesSpec (Proxy :: Proxy UnionType )
1723   (Proxy :: Proxy Value ) True
1724 ,typesSpec (Proxy :: Proxy (Counted NumberConstraint))
1725   (Proxy :: Proxy Scientific ) False
1726 ]
1727 representationSpec
1728
1729 typesSpec :: (Typeable ty
1730   ,Typeable term
1731   ,Monoid ty
1732   ,Arbitrary ty
1733   ,Arbitrary term
1734   ,Show ty
1735   ,Show term
1736   ,Eq ty
1737   ,Eq term
1738   ,Typelike ty
1739   ,Types ty term
1740   ,TypeCost ty
1741   )
1742   => Proxy ty
1743   -> Proxy term
1744   -> Bool -- idempotent?
1745   -> (String, [Laws])
1746 typesSpec (tyProxy :: Proxy ty)
1747   (termProxy :: Proxy term) isIdem =
1748   (nameOf @ty <> " types " <> nameOf @term, [
1749     arbitraryLaws tyProxy
1750   , eqLaws tyProxy
1751   , monoidLaws tyProxy
1752   , commutativeMonoidLaws tyProxy
1753   , typeCostLaws tyProxy
1754   , typelikeLaws tyProxy
1755   , arbitraryLaws termProxy
1756   , eqLaws termProxy
1757   , typesLaws tyProxy termProxy
1758   ]<>idem)
1759 where
1760 
```

```

1761     idem | isIdem    = [idempotentSemigroupLaws tyProxy]
1762         | otherwise = []
1763
1764 typesLaws :: (      ty `Types` term
1765             ,Arbitrary ty
1766             ,Arbitrary term
1767             ,Show    ty
1768             ,Show    term
1769             )
1770     => Proxy ty
1771     -> Proxy term
1772     -> Laws
1773 typesLaws (_ :: Proxy ty) (_ :: Proxy term) =
1774     Laws "Types" [("empty contains no terms"
1775                   ,property $
1776                     empty_contains_no_terms @ty @term)
1777                  ,("beyond contains all terms"
1778                   ,property $
1779                     beyond_contains_all_terms @ty @term)
1780                  ,("inferred type contains its term"
1781                   ,property $
1782                     inferred_type_contains_its_term @ty @term)
1783                  ]
1784
1785 <<representation-examples>>
1786
1787 representationTest :: String -> [Value] -> HType -> IO Bool
1788 representationTest name values repr = do
1789     if foundRepr == repr
1790     then do
1791         putStrLn $ "**** Representation test " <> name <> " succeeded."
1792         return True
1793     else do
1794         putStrLn $ "**** Representation test " <> name <> " failed."
1795         putStrLn $ "Values      : " <> show values
1796         putStrLn $ "Inferred type : " <> show inferredType
1797         putStrLn $ "Representation: " <> show foundRepr
1798         putStrLn $ "Expected     : " <> show repr
1799         return False
1800
1801 where
1802     foundRepr :: HType
1803     foundRepr = toHType inferredType
1804     inferredType :: UnionType
1805     inferredType = foldMap infer values
1806
1807 readJSON :: HasCallStack
1808     => BS.ByteString -> Value
1809 readJSON = fromMaybe ("Error reading JSON file")
1810     . decodeStrict
1811     . BS.unlines
1812     . filter notComment
1813     . BS.lines
1814
1815     notComment (BS.isPrefixOf "//" -> True) = False
1816     notComment _ = True
1817
1818 representationSpec :: IO ()
1819 representationSpec = do
1820     b <- sequence
1821     [representationTest "1a" example1a_values example1a_repr
1822     ,representationTest "1b" example1b_values example1b_repr
1823     ,representationTest "1c" example1c_values example1c_repr
1824     ,representationTest "2" example2_values example2_repr
1825     ,representationTest "3" example3_values example3_repr
1826     ,representationTest "4" example4_values example4_repr
1827     ,representationTest "5" example5_values example5_repr
1828     ,representationTest "6" example6_values example6_repr]
1829     when (not $ and b) $
1830         exitFailure
1831
1832
1833 Appendix: package dependencies
1834
1835 name: union-types
1836 version: '0.1.0.0'
1837 category: Web
1838 author: Anonymous
1839 maintainer: example@example.com
1840 license: BSD-3
1841 extra-source-files:
1842     CHANGELOG.md
1843     - README.md
1844 dependencies:
1845     - base
1846     - aeson
1847     - containers
1848     - text
1849     - hspectest
1850     - QuickCheck
1851     - unordered-containers
1852     - scientific
1853     - hspectest
1854     - QuickCheck
1855     - validity
1856     - vector
1857     - unordered-containers
1858     - scientific
1859     - genvalidity
1860     - genvalidity-hspec
1861     - genvalidity-property
1862     - time
1863     - email-validate
1864     - generic-arbitrary
1865     - mtl
1866     - hashable
1867 library:
1868     source-dirs: src
1869     exposed-modules:
1870

```

```

1871 - Unions
1872 tests:
1873 spec:
1874   main: Spec.hs
1875   source-dirs:
1876     - test/lib
1877     - test/spec
1878   dependencies:
1879     - union-types
1880     - mtl
1881     - random
1882     - transformers
1883     - hashable
1884     - quickcheck-classes
1885     - file-embed
1886     - bytestring
1887 less-arbitrary:
1888   main: LessArbitrary.hs
1889   source-dirs:
1890     - test/lib
1891     - test/less
1892   dependencies:
1893     - union-types
1894     - mtl
1895     - random
1896     - transformers
1897     - hashable
1898     - quickcheck-classes
1899     - quickcheck-instances

```

Appendix: representation of generated Haskell types

We will not delve here into identifier conversion between JSON and Haskell, so it suffices that we have an abstract datatypes for Haskell type and constructor identifiers:

```

1907 newtype HConsId = HConsId String
1908   deriving (Eq, Ord, Show, Generic, IsString)
1909 newtype HFieldId = HFieldId String
1910   deriving (Eq, Ord, Show, Generic, IsString)
1911 newtype HTypeId = HTypeId String
1912   deriving (Eq, Ord, Show, Generic, IsString)

```

For each single type we will either describe its exact representation or reference to the other definition by name:

```

1915 data HType =
1916   HRef HTypeId
1917 | HApp HTypeId [HType]
1918 | HADT [HCons]
1919   deriving (Eq, Ord, Show, Generic)

```

For syntactic convenience, we will allow string literals to denote type references:

```

1923 instance IsString HType where
1924   fromString = HRef . fromString

```

When we define a single constructor, we allow field and constructor names to be empty strings (""), assuming that the relevant identifiers will be put there by post-processing that will pick names using types of fields and their containers[???].

```

1931 data HCons = HCons {
1932     name :: HConsId
1933     , args :: [(HFieldId, HType)]
1934   }
1935   deriving (Eq, Ord, Show, Generic)

```

At some stage we want to split representation into individually named declarations, and then we use environment of defined types, with an explicitly named toplevel type:

```

1940 data HTypeEnv = HTypeEnv {
1941     toplevel :: HTypeId
1942     , env     :: HashMap HTypeId HType
1943   }

```

When checking for validity of types and type environments, we might need a list of predefined identifiers that are imported:

```

1947 predefinedHTypes :: [HType]
1948 predefinedHTypes = [
1949     "Data.Aeson.Value"
1950   , "()"
1951   , "Double"
1952   , "String"
1953   , "Int"
1954   , "Date" -- actually: "Data.Time.CalendarDay"
1955   , "Email" -- actually: "Data.Email"
1956 ]

```

Consider that we also have an htop value that represents any possible JSON value. It is polymorphic for ease of use:

```

1960 htop :: IsString s => s
1961 htop = "Data.Aeson.Value"

```

5.4 Code for selecting representation

Below is the code to select representation, as described in sec. 4.

To convert union type discipline to strict Haskell type representations, we need to join the options to get the actual representation:

```

1969 toHType :: ToHType ty => ty -> HType
1970 toHType = joinAlts . toHTypes

```

```

1972 joinAlts :: [HType] -> HType
1973 joinAlts [] = htop -- promotion of empty type
1974 joinAlts alts = foldr1 joinPair alts
1975   where
1976     joinPair a b = HApp "[:]" [a, b]

```

Considering the assembly of UnionType, we join all the options, and convert nullable types to Maybe types


```

1981 instance ToHType UnionType where
1982   toHTypes UnionType {..} =
1983     prependNullable unionNull opts
1984   where
1985     opts = concat [toHTypes unionBool
1986                   ,toHTypes unionStr
1987                   ,toHTypes unionNum
1988                   ,toHTypes unionArr
1989                   ,toHTypes unionObj]
1990
1991 prependNullable :: PresenceConstraint a -> [HType] -> [HType]
1992 prependNullable Present tys = [HApp "Maybe" [joinAlts tys]]
1993 prependNullable Absent tys = tys
1994
1995 The type class returns a list of mutually exclusive type
1996 representations:
1997 class Typelike ty
1998   => ToHType ty where
1999   toHTypes :: ty -> [HType]
2000
2001 Conversion of flat types is quite straightforward:
2002 instance ToHType BoolConstraint where
2003   toHTypes Absent = []
2004   toHTypes Present = ["Bool"]
2005 instance ToHType NumberConstraint where
2006   toHTypes NCNever = []
2007   toHTypes NCFloat = ["Double"]
2008   toHTypes NCInt = ["Int"]
2009 instance ToHType StringConstraint where
2010   toHTypes SCAny = ["String"]
2011   toHTypes SCEmail = ["Email"]
2012   toHTypes SCDate = ["Date"]
2013   toHTypes (SCEnum es) = [HADT $
2014     mkCons <$> Set.toList es
2015   ]
2016   where
2017     mkCons = (`HCons` [])
2018     . HConsId
2019     . Text.unpack
2020   toHTypes SCNever = []
2021
2022 For array and object types we pick the representation
2023 which presents the lowest cost of optionality:
2024 instance ToHType ObjectConstraint where
2025   toHTypes ObjectNever = []
2026   toHTypes ObjectConstraint {..} =
2027     if typeCost recordCase <= typeCost mappingCase
2028     then toHTypes recordCase
2029     else toHTypes mappingCase
2030 instance ToHType RecordConstraint where
2031   toHTypes RCBottom = []
2032   toHTypes RCTop = [htop] -- should never happen
2033   toHTypes (RecordConstraint fields) =
2034     [HADT
2035       [HCons "" $ fmap convert $ Map.toList fields]
2036     ]
2037   where
2038     convert (k,v) = (HFieldId $ Text.unpack k
2039       ,toHType v)
2040
2041 instance ToHType MappingConstraint where
2042   toHTypes MappingNever = []
2043   toHTypes MappingConstraint {..} =
2044     [HApp "Map" [toHType keyConstraint
2045       ,toHType valueConstraint
2046     ]]
2047
2048 instance ToHType RowConstraint where
2049   toHTypes RowNever = []
2050   toHTypes RowTop = [htop]
2051   toHTypes (Row cols) =
2052     [HADT
2053       [HCons "" $ fmap (\ut -> ("", toHType ut)) cols]
2054     ]
2055
2056 instance ToHType ArrayConstraint where
2057   toHTypes ArrayNever = []
2058   toHTypes ArrayConstraint {..} =
2059     if typeCost arrayCase <= typeCost rowCase
2060     -- || count <= 3
2061     then [toHType arrayCase]
2062     else [toHType rowCase ]
2063
2064
2065 Appendix: Missing pieces of code
2066
2067 In order to represent FreeType for the Value, we need to
2068 add Ord instance for it:
2069 instance Ord Value where
2070   compare = compare `on` hash
2071
2072 For validation of dates and emails, we import functions
2073 from Hackage:
2074 isValidDate :: Text -> Bool
2075 isValidDate = isJust
2076   . parseDate
2077   . Text.unpack
2078
2079 where
2080   parseDate :: String -> Maybe Day
2081   parseDate = parseTimeM True
2082     defaultTimeLocale $
2083     iso8601DateFormat Nothing
2084
2085 isValidEmail :: Text -> Bool
2086 isValidEmail = Text.Email.Validate.isValid
2087   . Text.encodeUtf8
2088
2089
2090 Appendix: Damas-Milner as Typelike

```