

GADTs meet subtyping

Gabriel Scherer and Didier Rémy

INRIA, Rocquencourt*

Abstract. While generalized algebraic datatypes (GADTs) are now considered well-understood, adding them to a language with a notion of subtyping comes with a few surprises. What does it mean for a GADT parameter to be covariant? The answer turns out to be quite subtle. It involves fine-grained properties of the subtyping relation that raise interesting design questions. We allow variance annotations in GADT definitions, study their soundness, and present a sound and complete algorithm to check them. Our work may be applied to real-world ML-like languages with explicit subtyping such as OCaml, or to languages with general subtyping constraints.

Introduction

In languages that have a notion of subtyping, the interface of parametrized types usually specifies a *variance*. It defines the subtyping relation between two instances of a parametrized type from the subtyping relations that hold between their parameters. For example, the type α `list` of immutable lists is expected to be *covariant*: we wish σ `list` \leq σ' `list` as soon as $\sigma \leq \sigma'$.

Variance is essential in languages with parametric polymorphism whose programming idioms rely on subtyping, in particular object-oriented languages, or languages with structural datatypes such as extensible records and variants, dependently typed languages with inductive types (to represent positivity requirements), or additional information in types such as permissions, effects, *etc.* A last reason to care about variance is its use in the *relaxed value restriction* [Gar04]: while a possibly-effectful expression, also called an *expansive expression*, cannot be soundly generalized in ML—unless some sophisticated enhancement of the type system keeps track of effectful expressions—it is always sound to generalize type variables that only appear in covariant positions, as they may not classify mutable data. Therefore, it is important for extensions of type definitions, such as generalized algebraic datatypes (GADTs), to support it as well through a clear and expressive definition of parameter covariance.

For example, consider the following GADT of well-typed expressions:

```
type + $\alpha$  exp =  
  | Val :  $\alpha \rightarrow \alpha$  exp  
  | Int : int  $\rightarrow$  int exp  
  | Thunk :  $\forall \beta. \beta$  exp * ( $\beta \rightarrow \alpha$ )  $\rightarrow$   $\alpha$  exp  
  | Prod :  $\forall \beta \gamma. \beta$  exp *  $\gamma$  exp  $\rightarrow$  ( $\beta * \gamma$ ) exp
```

* Part of this work has been done at IRILL.

Is it safe to say that `exp` is covariant in its type parameter? It turns out that, using the subtyping relation of the OCaml type system, the answer is “yes”. But, surprisingly to us, in a type system with a top type \top , the answer would be “no”. We introduce this example in details in §1—and present some interesting counter-examples of incorrect variance annotations.

Verifying variance annotations for simple algebraic datatypes is straightforward: it suffices to check that covariant type variables appear only positively and contravariant variables only negatively in the types of the arguments of the datatype constructors. GADTs can be formalized as extensions of datatypes where constructors have typed arguments, but also a set of existential variables and equality constraints. Then, the simple check of algebraic datatypes apparently becomes a searching problem: witnesses for existentials must be found so as to satisfy the equality constraints. That is, there is a natural correctness criterion (already present in previous work); however, it is expressed in a “semantic” form that is not suitable for a simple implementation in a type checker. We present this semantic criterion in §2 after reviewing the formal framework of variance-based subtyping.

The main contribution of our work, described in §3, is to develop a syntactic criterion that ensures the semantics criterion. Our solution extends the simple check of algebraic datatypes in a non-obvious way by introducing two new notions. First, *upward and downward-closure* of type constructors explains how to check that a single equality constraint is still satisfiable in presence of variance (but also raises interesting design issues for the subtyping relation). Second, *zipping* explains when witnesses exist for existential variables, that is, when multiple constraints using the same existential may soundly be used without interfering with each other. These two properties are combined into a new syntactic judgment of *decomposability* that is central to our syntactic criterion. We prove that our syntactic criterion is sound and complete with respect to the semantic criterion. The proof of soundness is relatively direct, but completeness is much harder.

We discuss the implication of our results in §4, in particular the notion of upward and downward-closure properties of type constructors, on the design of a subtyping relation. We also contrast this approach, motivated by the needs of a language of a ML family, with a different and mostly orthogonal approach taken by existing object-oriented languages, namely $C\sharp$ and Scala, where a natural notion of GADTs involves subtyping constraints, rather than equality constraints. We can re-evaluate our syntactic criterion in this setting: it is still sound, but the question of completeness is left open.

In summary, we propose a syntactic criterion for checking the soundness of variance annotations of GADTs with equality constraints in a language with subtyping. Our work is directly applicable to the OCaml language, but our approach can also be transposed to languages with general subtyping constraints, and raises interesting design questions. A long version of the present article, containing the detailed proofs and additional details and discussion, is available online [SR].

1 Examples

Let us first explain why it is reasonable to say that $\alpha \text{ exp}$ is covariant. Informally, if we are able to coerce a value of type α into one of type α' (we write $(v :> \alpha')$ to explicitly cast a value v of type α to a value of type α'), then we are also able to transform a value of type $\alpha \text{ exp}$ into one of type $\alpha' \text{ exp}$. Here is some pseudo-code¹ for the coercion function:

```
let coerce :  $\alpha \text{ exp} \rightarrow \alpha' \text{ exp}$  = function
| Val (v :  $\alpha$ ) -> Val (v :>  $\alpha'$ )
| Int n -> Int n
| Thunk  $\beta$  (b :  $\beta \text{ exp}$ ) (f :  $\beta \rightarrow \alpha$ ) ->
  Thunk  $\beta$  b (fun x -> (f x :>  $\alpha'$ ))
| Prod  $\beta$   $\gamma$  ((b, c) :  $\beta \text{ exp} * \gamma \text{ exp}$ ) ->
  (* if  $\beta * \gamma \leq \alpha'$ , then  $\alpha'$  is of the form  $\beta' * \gamma'$ 
    with  $\beta \leq \beta'$  and  $\gamma \leq \gamma'$  *)
  Prod  $\beta'$   $\gamma'$  ((b :>  $\beta' \text{ exp}$ ), (c :>  $\gamma' \text{ exp}$ ))
```

In the `Prod` case, we make an informal use of something we know about the OCaml type system: the supertypes of a tuple are all tuples. By entering the branch, we gain the knowledge that α must be equal to some type of the form $\beta * \gamma$. So from $\alpha \leq \alpha'$ we know that $\beta * \gamma \leq \alpha'$. Therefore, α' must itself be a pair of the form $\beta' * \gamma'$. By covariance of the product, we deduce that $\beta \leq \beta'$ and $\gamma \leq \gamma'$. We may thus conclude by casting at types $\beta' \text{ exp}$ and $\gamma' \text{ exp}$, recursively.

Similarly, in the `Int` case, we know that α must be an `int` and therefore an `int exp` is returned. This is because we know that, in OCaml, no type is above `int`: if $\text{int} \leq \tau$, then τ must be `int`.

What we use in both cases is reasoning of the form²: “if $T[\bar{\beta}] \leq \alpha'$, then I know that α' is of the form $T[\bar{\beta}']$ for some $\bar{\beta}'$ ”. We call this an *upward closure* property: when we “go up” from a $T[\bar{\beta}]$, we only find types that also have the structure of T . Similarly, for contravariant parameters, we would need a *downward closure* property: T is downward-closed if $T[\bar{\beta}] \geq \alpha'$ entails that α' is of the form $T[\bar{\beta}']$.

Before studying a more troubling example, we define the classic equality type $(\alpha, \beta) \text{ eq}$ and the corresponding casting function $\text{cast} : \forall \alpha \beta. (\alpha, \beta) \text{ eq} \rightarrow \alpha \rightarrow \beta$:

```
type ( $\alpha$ ,  $\beta$ ) eq = let cast r =
| Refl :  $\forall \gamma. (\gamma, \gamma) \text{ eq}$  match r with Refl -> (fun x -> x)
```

Notice that it would be unsound³ to define `eq` as covariant, even in only one parameter. For example, if we had `type (+ α , = β) eq`, from any $\sigma \leq \tau$, we could subtype $(\sigma, \sigma) \text{ eq}$ into $(\tau, \sigma) \text{ eq}$, allowing a cast from any value of type τ back into one of type σ , which is unsound in general.

¹ The variables β' and γ' of the `Prod` case are never really defined, only justified at the meta-level, making this code only an informal sketch.

² We write $T[\bar{\beta}]$ for a type expression T that may contain free occurrences of variables $\bar{\beta}$ and $T[\bar{\sigma}]$ for the simultaneous substitution of $\bar{\sigma}$ for $\bar{\beta}$ in T .

³ This counterexample is due to Jeremy Yallop.

As a counter-example, the following declaration is incorrect: the type α `bad` cannot be declared covariant.

```
type + $\alpha$  bad =
  | K : < m : int >  $\rightarrow$  < m : int > bad
let v = (K (object method m = 1 end) )> < > bad)
```

This declaration uses the OCaml object type `< m : int >`, which qualifies objects having a method `m` returning an integer. It is a subtype of object types with fewer methods, in this case the empty object type `< >`, so the alleged covariance of `bad`, if accepted by the compiler, would allow us to cast a value of type `< m : int > bad` into one of type `< > bad` and thus have the above value `v` of type `< > bad`. However, if such a value `v` existed, we could produce an equality witness `(< >, < m : int >) eq` that allows to cast any empty object of type `< >` into an object of type `< m : int >`, but this is unsound, of course!

```
let get_eq :  $\alpha$  bad  $\rightarrow$  ( $\alpha$ , < m : int >) eq = function
  | K _ -> Refl      (* locally  $\alpha = < m : int > *$ )
let wrong : < > -> < m : int > =
  let eq : (< >, < m : int >) eq = get_eq v in cast eq
```

It is possible to reproduce this example using a different feature of the OCaml type system named *private type abbreviation*⁴: a module using a type `type s = τ internally` may describe its interface as `type s = private τ` . This is a compromise between a type abbreviation and an abstract type: it is possible to cast a value of type `s` into one of type `τ` , but not, conversely, to construct a value of type `s` from one of type `τ` . In other words, `s` is a strict subtype of `τ` : we have `s \leq τ` but not `s \geq τ` . Take for example `type file_descr = private int`: this semi-abstraction is useful to enforce invariants by restricting the construction of values of type `file_descr`, while allowing users to conveniently and efficiently destruct them for inspection at type `int`. Using an unsound but quite innocent-looking covariant GADT datatype, one is able to construct a function to cast any integer into a `file_descr`, which defeats the purpose of this abstraction—see the extended version of this article for the full example.

The difference between the former, correct `Prod` case and those two latter situations with unsound variance is the notion of upward closure. The types `$\alpha * \beta$` and `int` used in the correct example were upward-closed. On the contrary, the private type `file_descr` has a distinct supertype `int`, and similarly, the object type `< m:int >` has a supertype `< >` with a different structure (no method `m`).

Finally, the need for covariance of `α exp` can be justified either by applications using subtyping on data (for example object types or polymorphic variants), or by the relaxed value restriction. If we used the `Thunk` constructor to delay a computation returning an object of type `< m : int >`, that is itself of type `< m : int > exp`, we may need to see it as a computation returning the empty object `< >`. We could also wish to define an abstract interface through a module boundary that would not expose any implementation detail about the datatype; for example, using `Product` to implement a `list` interface.

⁴ This counterexample is due to Jacques Garrigue.

```

module Exp : sig
  type  $\alpha$  exp
  val inj :  $\alpha \rightarrow \alpha$  exp
  val pair :  $\alpha$  exp  $\rightarrow$   $\beta$  exp  $\rightarrow$  ( $\alpha * \beta$ ) exp
  val fst : ( $\alpha * \beta$ ) exp  $\rightarrow$   $\alpha$  exp
end

```

What would then be the type of `Exp.inj []`? In presence of the value restriction, this application cannot be generalized, and we get a weak polymorphic type `? α list Exp.exp` for some non-generalized inference variable `? α` . If we change the interface to express that `Exp.exp` is covariant, then we get the expected polymorphic type `$\forall \alpha. \alpha$ list Exp.exp`.

2 A formal setting

2.1 The subtyping relation

Ground types consist of base type `q`, types `τ p`, function types `$\tau_1 \rightarrow \tau_2$` , product types `$\tau_1 * \tau_2$` , and a set of algebraic datatypes `$\bar{\sigma}$ t`. We also write `σ` and `ρ` for types, `$\bar{\sigma}$` for a sequence of types `$(\sigma_i)_{i \in I}$` , and we use prefix notation for datatype parameters, as is the usage in ML. Datatypes may be user-defined by toplevel declarations of the form:

$$\text{type } \bar{v}\bar{\alpha} \text{ t} = K_1 \text{ of } \tau^1[\bar{\alpha}] \mid \dots \mid K_n \text{ of } \tau^n[\bar{\alpha}]$$

This is a disjoint sum: the constructors `K_c` represent all possible cases and each type `$\tau^c[\bar{\alpha}]$` is the domain of the constructor `K_c` . Applying `K_c` to an argument `e` of a corresponding ground type `$\tau[\bar{\sigma}]$` constructs a term of type `$\bar{\sigma} \text{ t}$` . Values of this type are deconstructed using pattern matching clauses of the form `$K_c x \rightarrow e$` , one for each constructor.

The sequence `$\bar{v}\bar{\alpha}$` is a binding list of type variables `α_i` along with their *variance annotation* `v_i` . Variances range in the set `$\{+, -, =, \bowtie\}$` . We may associate a relation `\prec_v` between types to each variance `v` :

- `\prec_+` is the *covariant* relation (`\leq`);
- `\prec_-` is the *contravariant* relation (`\geq`), the symmetric of (`\leq`);
- `$\prec_ =$` is the *invariant* relation (`$=$`) defined as the intersection of (`\leq`) and (`\geq`);
- `\prec_{\bowtie}` is the *irrelevant* relation (`\bowtie`), *i.e.* the full relation such that `$\sigma \bowtie \tau$` holds for all types `σ` and `τ` .

Given a reflexive transitive relation (`\leq`) on base types, the subtyping relation on ground types (`\leq`) is defined by the inference rules of Figure 1, which, in particular, give their meaning to the variance annotations `$\bar{v}\bar{\alpha}$` . The judgment `type $\bar{v}\bar{\alpha} \text{ t}$` simply means that the type constructor `t` has been previously defined with the variance annotation `$\bar{v}\bar{\alpha}$` . Notice that the rules for arrow and product types, `SUB-FUN` and `SUB-PROD`, can be subsumed by the rule for datatypes `SUB-CONSTR`. Indeed, one can consider them as special datatypes (with a specific

| | | | | |
|--------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------|--|
| $\frac{\text{SUB-REFL}}{\sigma \leq \sigma}$ | $\frac{\text{SUB-TRANS} \quad \sigma_1 \leq \sigma_2 \quad \sigma_2 \leq \sigma_3}{\sigma_1 \leq \sigma_3}$ | $\frac{\text{SUB-FUN} \quad \sigma \geq \sigma' \quad \tau \leq \tau'}{\sigma \rightarrow \tau \leq \sigma' \rightarrow \tau'}$ | | |
| $\frac{\text{SUB-PROD} \quad \sigma \leq \sigma' \quad \tau \leq \tau'}{\sigma * \tau \leq \sigma' * \tau'}$ | $\frac{\text{SUB-CONSTR} \quad \text{type } \bar{v}\bar{\alpha} \mathbf{t} \quad \forall i, \sigma_i \prec_{v_i} \sigma'_i}{\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}}$ | $\frac{\text{SUB-P} \quad \sigma \leq \sigma'}{\sigma \mathbf{p} \leq \sigma' \mathbf{p}}$ | $\frac{\text{SUB-PQ}}{\sigma \mathbf{p} \leq \mathbf{q}}$ | |

Fig. 1. Subtyping relation

dynamic semantics) of variance $(-, +)$ and $(+, +)$, respectively. For this reason, the following definitions will not explicitly detail the cases for arrows and products.

The rules SUB-P and SUB-PQ were added for the explicit purpose of introducing some amount of non-atomic subtyping in our relation. For two fixed type constructors \mathbf{p} (unary) and \mathbf{q} (nullary), we have $\sigma \mathbf{p} \leq \mathbf{q}$ for any σ . Note that \mathbf{q} is not a top type as it is not above all types, only above the $\sigma \mathbf{p}$. Of course, we could add other such type constructors, but those are enough to make the system interesting and representative of complex subtype relation.

As usual in subtyping systems, we could reformulate our judgment in a syntax-directed way, to prove that it admits good inversion properties: if $\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}$ and $\text{type } \bar{v}\bar{\alpha} \mathbf{t}$, then one can deduce that for each i , $\sigma_i \prec_{v_i} \sigma'_i$.

The non-atomic rule SUB-PQ ensures that our subtyping relation is not “too structured” and is a meaningful choice for a formal study applicable to real-world languages with possibly top or bottom types, private types, record with subtyping, etc. In particular, the type constructor \mathbf{p} is *not* upward-closed (and conversely \mathbf{q} is not downward-closed), as used informally in the examples and defined for arbitrary variances in the following way:

Definition 1 (Constructor closure). *A type constructor $\bar{\alpha} \mathbf{t}$ is v -closed if, for any type sequence $\bar{\sigma}$ and type τ such that $\bar{\sigma} \mathbf{t} \prec_v \tau$ hold, then τ is necessarily equal to $\bar{\sigma}' \mathbf{t}$ for some $\bar{\sigma}'$.*

2.2 The algebra of variances

If we know that $\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}$, that is $\bar{\sigma} \mathbf{t} \prec_+ \bar{\sigma}' \mathbf{t}$, and the constructor \mathbf{t} has variable $\bar{v}\bar{\alpha}$, an inversion principle tells us that $\sigma_i \prec_{v_i} \sigma'_i$ for each i . But what if we only know $\bar{\sigma} \mathbf{t} \prec_u \bar{\sigma}' \mathbf{t}$ for some variance u different from $(+)$? If u is $(-)$, we get the reverse relation $\sigma_i \succ_{v_i} \sigma'_i$. If u is (\bowtie) , we get $\sigma_i \bowtie \sigma'_i$, that is, nothing. This outlines a *composition* operation on variances $u.v_i$, such that if $\bar{\sigma} \mathbf{t} \prec_u \bar{\sigma}' \mathbf{t}$ then $\sigma_i \prec_{u.v_i} \sigma'_i$ holds. It is defined by the table in figure 2.2.

This operation is associative and commutative. Such an operator, and the algebraic properties of variances explained below, have already been used by other authors, for example [Abe06].

There is a natural order relation between *variances*, which is the *coarser-than* order between the corresponding relations: $v \leq w$ if and only if $(\prec_v) \supseteq (\prec_w)$;

i.e. if and only if, for all σ and τ , $\sigma \prec_w \tau$ implies $\sigma \prec_v \tau$.⁵ This reflexive, partial order is described by the lattice diagram in figure 2.2. All variances are smaller than = and bigger than \bowtie .

| $v.w$ | = | + | - | \bowtie | w |
|-----------|-----------|-----------|-----------|-----------|-----------|
| = | = | = | = | \bowtie | \bowtie |
| + | = | + | - | \bowtie | \bowtie |
| - | = | - | + | \bowtie | \bowtie |
| \bowtie | \bowtie | \bowtie | \bowtie | \bowtie | \bowtie |
| v | | | | | |

Fig. 2. Variance composition table

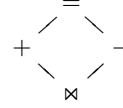


Fig. 3. Variance order diagram

From the order lattice on variances we can define join \vee and meet \wedge of variances: $v \vee w$ is the biggest variance such that $v \vee w \leq v$ and $v \vee w \leq w$; conversely, $v \wedge w$ is the lowest variance such that $v \leq v \wedge w$ and $w \leq v \wedge w$. Finally, the composition operation is monotone: if $v \leq v'$ then $w.v \leq w.v'$ and $v.w \leq v'.w$.

We often manipulate vectors $\overline{v\alpha}$ of variable associated with variances, which correspond to the “context” Γ of a type declaration. We extend our operation pairwise on those contexts: $\Gamma \vee \Gamma'$ and $\Gamma \wedge \Gamma'$, and the ordering between contexts $\Gamma \leq \Gamma'$. We also extend the variance-dependent subtyping relation (\prec_v), which becomes an order (\prec_Γ) between vectors of type of the same length: $\overline{\sigma} \prec_{\overline{v\alpha}} \overline{\sigma}'$ holds when we have $\sigma_i \prec_{v_i} \sigma'_i$ for all i .

2.3 A judgment for variance of type expressions

We define a judgment to check the variance of a type expression. Given a context Γ of the form $\overline{v\alpha}$, that is, where each variable is annotated with a variance, the judgment $\Gamma \vdash \tau : v$ checks that the expression τ varies along v when the variables of τ vary along their variance in Γ . For example, $(+\alpha) \vdash \tau[\alpha] : +$ holds when $\tau[\alpha]$ is covariant in its variable α . The inference rules for the judgment $\Gamma \vdash \tau : v$ are defined on Figure 4.

The parameter v evolves when going into subderivations: when checking $\Gamma \vdash \tau_1 \rightarrow \tau_2 : v$, contravariance is expressed by checking $\Gamma \vdash \tau_1 : (v.-)$. Previous work (on variance as [Abe06] and [EKRY06], but also on irrelevance as in [Pfe01]) used no such parameter, but modified the context instead, checking $\Gamma / - \vdash \tau_1$ for some “variance cancellation” operation $vw/$ (see [Abe06] for a principled

⁵ The reason for this order reversal is that the relations occur as hypotheses, in negative position, in definition of subtyping: if we have $v \leq w$ and **type** $v\alpha \mathbf{t}$, it is safe to assume **type** $w\alpha \mathbf{t}$, since $\sigma \prec_w \sigma'$ implies $\sigma \prec_v \sigma'$, which implies $\sigma \mathbf{t} \leq \sigma' \mathbf{t}$. One may also see it, as Abel notes, as an “information order”: knowing that $\sigma \prec_+ \tau$ “gives you more information” than knowing that $\sigma \prec_{\bowtie} \tau$, therefore $\bowtie \leq +$.

$$\begin{array}{c}
\text{VC-VAR} \\
\frac{w\alpha \in \Gamma \quad w \geq v}{\Gamma \vdash \alpha : v} \\
\\
\text{VC-CONSTR} \\
\frac{\Gamma \vdash \text{type } \bar{w}\bar{\alpha} \mathbf{t} \quad \forall i, \Gamma \vdash \sigma_i : v.w_i}{\Gamma \vdash \bar{\sigma} \mathbf{t} : v}
\end{array}$$

Fig. 4. Variance assignment

presentation). Our own inference rules preserve the same context in the whole derivation and can be more easily adapted to the decomposability judgment $\Gamma \vdash \tau : v \Rightarrow v'$ that we introduce in §3.4.

A semantics for variance assignment This syntactic judgment $\Gamma \vdash \tau : v$ corresponds to a semantic property about the types and context involved, which formalizes our intuition of “when the variables vary along Γ , the expression τ varies along v ”. We also give a few formal results about this judgment.

Definition 2 (Interpretation of the variance checking judgment).

We write $\llbracket \Gamma \vdash \tau : v \rrbracket$ for the property: $\forall \bar{\sigma}, \bar{\sigma}', \bar{\sigma} \prec_{\Gamma} \bar{\sigma}' \implies \tau[\bar{\sigma}] \prec_v \tau[\bar{\sigma}']$.

Lemma 1 (Correctness of variance checking).

$\Gamma \vdash \tau : v$ is provable if and only if $\llbracket \Gamma \vdash \tau : v \rrbracket$ holds.

Lemma 2 (Monotonicity).

If $\Gamma \vdash \tau : v$ is provable and $\Gamma \leq \Gamma'$ then $\Gamma' \vdash \tau : v$ is provable.

Lemma 3 (Principality). For any type τ and any variance v , there exists a minimal context Δ such that $\Delta \vdash \tau : v$ holds. That is, for any other context Γ such that $\Gamma \vdash \tau : v$, we have $\Delta \leq \Gamma$.

We can generalize inversion of head type constructors (§2.1) to whole type expressions. The most general inversion is given by the principal context.

Theorem 1 (Inversion). For any type $\tau[\bar{\alpha}]$, variance v , and type sequences $\bar{\sigma}$ and $\bar{\sigma}'$, the subtyping relation $\tau[\bar{\sigma}] \prec_v \tau[\bar{\sigma}']$ holds if and only if the judgment $\Gamma \vdash \tau : v$ holds for some context Γ such that $\bar{\sigma} \prec_{\Gamma} \bar{\sigma}'$. Furthermore, if $\tau[\bar{\sigma}] \prec_v \tau[\bar{\sigma}']$ holds, then $\bar{\sigma} \prec_{\Delta} \bar{\sigma}'$ holds, where Δ is the minimal context such that $\Delta \vdash \tau : v$.

2.4 Variance annotations in ADTs

As a preparation for the difficult case of GADTs, we first present our approach in the well-understood case of algebraic datatypes. We exhibit a semantic criterion that justifies the correctness of a variance annotation; then, we propose an equivalent syntactic judgment. Of course, we recover the usual criterion that covariant variables should only occur positively.

In general, an ADT definition of the form

$$\text{type } \bar{w}\bar{\alpha} \mathbf{t} = \big|_{c \in C} K_c \text{ of } \tau^c[\bar{\alpha}]$$

cannot be accepted with any variance $\bar{v}\bar{\alpha} \mathbf{t}$. For example, the declaration (`type vα inv = Fun of α → α`) is only sound when v is invariant. Accepting a variance assignment $\bar{v}\bar{\alpha}$ determines the relations between closed types $\bar{\sigma}$ and $\bar{\sigma}'$ under which the relation $\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}$ is correct.

In the definition of $+\alpha \mathbf{exp}$ we justified the covariance of \mathbf{exp} by the existence of a coercion function. We now formalize this idea for the general case. To check the correctness of $\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}$ we check the existence of a coercion term that turns a closed value q of type $\bar{\sigma} \mathbf{t}$ into one of type $\bar{\sigma}' \mathbf{t}$ that is equal to q up to type information. We actually search for coercions of the form:

$$\mathbf{match} (q : \bar{\sigma} \mathbf{t}) \mathbf{with} \mid_{c \in C} K_c(x : \tau^c[\bar{\sigma}]) \rightarrow K_c(x :> \tau^c[\bar{\sigma}'])$$

Note that erasing types gives an η -expansion of the sum type, *i.e.* this is really a coercion. Hence, such a coercion exists if and only if it is well-typed, that is, each cast of the form $(x : \tau^c[\bar{\sigma}] :> \tau^c[\bar{\sigma}'])$ is itself well-typed. This gives our semantic criterion for ADTs.

Definition 3 (Semantic soundness criterion for ADTs).

We accept the ADT definition of $\bar{v}\bar{\alpha} \mathbf{t}$ with constructors $(K_c \text{ of } \tau^c[\bar{\alpha}])_{c \in C}$ if

$$\forall c \in C, \forall \bar{\sigma}, \forall \bar{\sigma}', \quad \bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t} \implies \tau^c[\bar{\sigma}] \leq \tau^c[\bar{\sigma}']$$

The syntactic criterion for ADTs We notice that this criterion is exactly the semantic interpretation of the variance checking judgment (Definition 2): the type `type $\bar{v}\bar{\alpha} \mathbf{t}$` is accepted if and only if the judgment $\bar{v}\bar{\alpha} \vdash \tau^c : (+)$ is derivable for each constructor type $\tau^c[\bar{\alpha}]$.

This syntactic criterion coincides with the well-known algorithm implemented in type checkers⁶: checking positive occurrences of a variable α corresponds to a proof obligation of the form $\bar{v}\bar{\alpha} \vdash \alpha : +$, which is valid only when α has variance $(+)$ or $(=)$ in Γ ; checking negative occurrences correspond to a proof obligation $\bar{v}\bar{\alpha} \vdash \alpha : -$, *etc.* This extends seamlessly to irrelevant variables, which must appear only under irrelevant context $\bar{v}\bar{\alpha} \vdash \alpha : \times$ —or not at all.

2.5 Variance annotations in GADTs

A general description of GADTs When used to build terms of type $\bar{\alpha} \mathbf{t}$, a constructor `K of τ` behaves like a function of type $\forall \bar{\alpha}. (\tau \rightarrow \bar{\alpha} \mathbf{t})$. Notice that the codomain is exactly $\bar{\alpha} \mathbf{t}$, the type \mathbf{t} instantiated with parametric variables. GADTs arise by relaxing this restriction, allowing constructors with richer types of the form $\forall \bar{\alpha}. (\tau \rightarrow \bar{\sigma} \mathbf{t})$. See for example the declaration of constructor `Prod` in the introduction:

$$\mid \mathbf{Prod} : \forall \beta \gamma. \beta \mathbf{exp} * \gamma \mathbf{exp} \rightarrow (\beta * \gamma) \mathbf{exp}$$

⁶ One should keep in mind that this criterion suffers the usual bane of static typing, it can reject programs that do not go wrong: `type $-\alpha \mathbf{weird} = K \text{ of } \alpha * \perp$` . For more details, see the beginning of the §3 in the long version of this article.

Instead of being just $\alpha \text{ exp}$, the codomain is now $(\beta * \gamma) \text{ exp}$. We moved from simple algebraic datatypes to so-called *generalized* algebraic datatypes. This approach is natural and convenient for the users, so it is exactly the syntax chosen in languages with explicit GADTs support, such as Haskell and OCaml, and is reminiscent of the inductive datatype definitions of dependently typed languages.

However, for the formal study of GADTs, a different formulation based on equality constraints is preferred. We use the following equivalent presentation, already present in previous works [SP07]. We force the codomain of the constructor Prod to be $\alpha \text{ t}$ again, instead of $(\beta * \gamma) \text{ t}$, by adding an explicit equality constraint $\alpha = \beta * \gamma$.

```

type  $\alpha \text{ exp} =$ 
  | Val of  $\exists \beta [\alpha = \beta]. \beta$ 
  | Int of  $[\alpha = \text{int}]. \text{int}$ 
  | Thunk of  $\exists \beta \gamma [\alpha = \gamma]. \beta \text{ exp} * (\beta \rightarrow \gamma)$ 
  | Prod of  $\exists \beta \gamma [\alpha = \beta * \gamma]. \beta \text{ exp} * \gamma \text{ exp}$ 

```

In the rest of the paper, we extend our former core language with such definitions. This does not impact the notion of subtyping, which is defined on GADT type constructors with variance $\text{type } \bar{v}\bar{\alpha} \text{ t}$ just as it previously was on simple ADT type constructors. What needs to be changed, however, is the soundness criterion for checking the variance of type definitions

The correctness criterion We must adapt our semantic criterion for datatype declarations (Definition 3) from simple ADTs to GADTs. Again, we check under which relations between $\bar{\sigma}$ and $\bar{\sigma}'$ the subtyping relation $\bar{\sigma} \text{ t} \leq \bar{\sigma}' \text{ t}$ holds for some GADT definition $\bar{v}\bar{\alpha} \text{ t}$.

The difference is that a constructor K_c that had an argument of type $\tau^c[\bar{\alpha}]$ in the simple ADT case, now has the more complex type $\exists \bar{\beta} [D[\bar{\alpha}, \bar{\beta}]]. \tau^c[\bar{\beta}]$, for a set of existential variables $\bar{\beta}$ and a set of equality constraints D —of the form $(\alpha_i = T_i[\bar{\beta}])_{i \in I}$ for a family of type expressions $(T_i[\bar{\beta}])_{i \in I}$. Given a closed value q of type $\bar{\sigma} \text{ t}$, the coercion term is:

$$\text{match } (q : \bar{\sigma} \text{ t}) \text{ with } |_{c \in C} K_c(x : \tau^c[\bar{\rho}_c]) \rightarrow K_c(x :> \tau^c[\bar{\rho}'_c])$$

We do not need to consider the dead cases: we only match on the constructors for which there exists an instantiation $\bar{\rho}_c$ of the existential variables $\bar{\beta}$ such that the constraint $D[\bar{\sigma}, \bar{\rho}]$, i.e. $\bigwedge_{i \in I} \sigma_i = T_i[\bar{\rho}_c]$, holds. To type-check this term, we need to find another instantiation $\bar{\rho}'_c$ that verifies the constraints $D[\bar{\sigma}', \bar{\rho}']$. This coercion type-checks only when $\tau^c[\bar{\rho}_c] \leq \tau^c[\bar{\rho}'_c]$ holds. This gives our semantic criterion for GADTs:

Definition 4 (Semantic soundness criterion for GADTs). *We accept the GADT definition of $\text{type } \bar{v}\bar{\alpha} \text{ t}$ with constructors $(K_c \text{ of } \exists \bar{\beta} [D[\bar{\alpha}, \bar{\beta}]]. \tau^c[\bar{\alpha}])_{c \in C}$, if for all c in C we have:*

$$\forall \bar{\sigma}, \bar{\sigma}', \bar{\rho}, (\bar{\sigma} \text{ t} \leq \bar{\sigma}' \text{ t} \wedge D[\bar{\sigma}, \bar{\rho}]) \implies \exists \bar{\rho}', D[\bar{\sigma}', \bar{\rho}'] \wedge \tau[\bar{\rho}] \leq \tau[\bar{\rho}'] \quad (\text{REQ})$$

As for ADTs, this criterion ensures soundness: if, under some variance annotation, a datatype declaration satisfies it, then the implied subtyping relations are all expressible as coercions in the language, and therefore correct. Whereas the simpler ADT criterion was already widely present in the literature, this one is less known; it is however present in the previous work of Simonet and Pottier [SP07] (presented as a constraint entailment problem).

Another way to understand this criterion would be to define constrained existential types of the form $\exists \bar{\beta}[D[\bar{\alpha}, \bar{\beta}]].\tau[\bar{\beta}]$ as first-class types and, with the right notion of subtyping for those, require that $\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}$ imply $(\exists \bar{\beta}[D[\bar{\sigma}, \bar{\beta}]].\tau[\bar{\beta}]) \leq (\exists \bar{\beta}[D[\bar{\sigma}', \bar{\beta}]].\tau[\bar{\beta}])$. The (easy) equivalence between those two presentations is detailed in the work of Simonet and Pottier [SP07].

3 Checking variances of GADT

3.1 Expressing decomposability

If we specialize REQ to the Prod constructor of the α exp example datatype, *i.e.* Prod of $\exists \beta \gamma[\alpha = \beta * \gamma]\beta \text{ exp} * \gamma \text{ exp}$, we get:

$$\forall \sigma, \sigma', \rho_1, \rho_2, \\ (\sigma \text{ exp} \leq \sigma' \text{ exp} \wedge \sigma = \rho_1 * \rho_2 \implies \exists \rho'_1, \rho'_2, (\sigma' = \rho'_1 * \rho'_2 \wedge \rho_1 * \rho_2 \leq \rho'_1 * \rho'_2))$$

We can substitute equalities and use the (user-defined) covariance to simplify the subtyping constraint $\sigma \text{ exp} \leq \sigma' \text{ exp}$ into $\sigma \leq \sigma'$:

$$\forall \sigma', \rho_1, \rho_2, (\rho_1 * \rho_2 \leq \sigma' \implies \exists \rho'_1, \rho'_2, (\sigma' = \rho'_1 * \rho'_2 \wedge \rho_1 \leq \rho'_1 \wedge \rho_2 \leq \rho'_2)) \quad (1)$$

This is the *upward closure* property mentioned in the introduction. The preceding transformation is safe only if any supertype σ' of a product $\rho_1 * \rho_2$ is itself a product, *i.e.* is of the form $\rho'_1 * \rho'_2$ for some ρ'_1 and ρ'_2 .

More generally, for a type $\Gamma \vdash \sigma$ and a variance v , we are interested in a closure property of the following form, where the notation $(\bar{\rho} : \Gamma)$ simply classifies type vectors $\bar{\rho}$ that have exactly one type ρ_i for each variable in Γ :

$$\forall (\bar{\rho} : \Gamma), \sigma', \quad \sigma[\bar{\rho}] \prec_v \sigma' \implies \exists (\bar{\rho}' : \Gamma), \sigma' = \sigma[\bar{\rho}']$$

Here, the context Γ represents the set of existential variables of the constructor (β and γ in our example). We can easily express the condition $\rho_1 \leq \rho'_1$ and $\rho_2 \leq \rho'_2$ on the right-hand side of the implication by considering a context Γ annotated with variances $(+\beta, +\gamma)$, and using the context ordering (\prec_Γ) . Then, (1) is equivalent to:

$$\forall (\bar{\rho} : \Gamma), \sigma', \quad \sigma[\bar{\rho}] \prec_v \sigma' \implies \exists (\bar{\rho}' : \Gamma), \bar{\rho} \prec_\Gamma \bar{\rho}' \wedge \sigma' = \sigma[\bar{\rho}']$$

Our aim is now to find a set of inference rules to check decomposability; we will later reconnect it to REQ. In fact, we study a slightly more general relation, where the equality $\sigma[\bar{\rho}'] = \sigma'$ on the right-hand side is relaxed to an arbitrary relation $\sigma[\bar{\rho}'] \prec_v \sigma'$:

Definition 5 (Decomposability). *Given a context Γ , a type expression $\sigma[\bar{\beta}]$ and two variances v and v' , we say that σ is decomposable under Γ from variance v to variance v' , which we write $\Gamma \Vdash \sigma : v \rightsquigarrow v'$, if the following property holds:*

$$\forall(\bar{\rho} : \Gamma), \sigma', \sigma[\bar{\rho}] \prec_v \sigma' \implies \exists(\bar{\rho}' : \Gamma), \bar{\rho} \prec_{\Gamma} \bar{\rho}' \wedge \sigma[\bar{\rho}'] \prec_{v'} \sigma'$$

We use the symbol \Vdash rather than \vdash to highlight the fact that this is just a logic formula, not the semantic interpretation of a syntactic judgment—we will introduce one later in section 3.4.

Remark that, due to the *positive* occurrence of the relation \prec_{Γ} in the proposition $\Gamma \Vdash \tau : v \rightsquigarrow v'$ and the anti-monotonicity of \prec_{Γ} , this formula is “anti-monotone” with respect to the context ordering $\Gamma \leq \Gamma'$. This corresponds to saying that we can still decompose, but with less information on the existential witness $\bar{\rho}'$.

Lemma 4 (Anti-monotonicity).

If $\Gamma \Vdash \tau : v \rightsquigarrow v'$ holds and $\Gamma' \leq \Gamma$, then $\Gamma' \Vdash \tau : v \rightsquigarrow v'$ also holds.

3.2 Variable occurrences

In the `Prod` case, the type whose decomposability was considered is $\beta * \gamma$ (in the context β, γ). In this very simple case, decomposability depends only on the type constructor for the product. In the present type system, with very strong invertibility principles on the subtyping relation, both upward and downward closures hold for products. In the general case, we require that this specific type constructor be upward-closed.

In general, the closure of the head type constructor alone is not enough to ensure decomposability of the whole type. For example, in a complex type expression with subterms, we should consider the closure of the type constructors appearing in the subterms as well. Besides, there are subtleties when a variable occurs several times.

For example, while $\beta * \gamma$ is decomposable from $(+)$ to $(=)$, $\beta * \beta$ is not: $\perp * \perp$ is an instantiation of $\beta * \beta$, and a subtype of, *e.g.*, `int * bool`, which is not an instance⁷ of $\beta * \beta$. The same variable occurring twice in covariant position (or having one covariant and one invariant or contravariant occurrence) breaks decomposability.

On the other hand, two invariant occurrences are possible: $\beta \text{ ref} * \beta \text{ ref}$ is upward-closed (assuming the type constructor `ref` is invariant and upward-closed): if $(\sigma \text{ ref} * \sigma \text{ ref}) \leq \sigma'$, then by upward closure of the product, σ' is of the form $\sigma'_1 * \sigma'_2$, and by its covariance $\sigma \text{ ref} \leq \sigma'_1$ and $\sigma \text{ ref} \leq \sigma'_2$. Now by invariance of `ref` we have $\sigma'_1 = \sigma \text{ ref} = \sigma'_2$, and therefore σ' is equal to $\sigma \text{ ref} * \sigma \text{ ref}$, which is an instance of $\beta \text{ ref} * \beta \text{ ref}$.

⁷ We use the term *instance* to denote the replacement of all the free variables of a type expression under context by closed types—not the specialization of an ML type scheme.

Finally, a variable may appear in irrelevant positions without affecting closure properties; $\beta * (\beta \text{ irr})$ (where irr is an upward-closed irrelevant type, defined for example as $\text{type } \alpha \text{ irr} = \text{int}$) is upward closed: if $\sigma * (\sigma \text{ irr}) \leq \sigma'$, then σ' is of the form $\sigma'_1 * (\sigma'_2 \text{ irr})$ with $\sigma \leq \sigma'_1$ and $\sigma \bowtie \sigma'_2$, which is equiconvertible to $\sigma'_1 * (\sigma'_1 \text{ irr})$ by irrelevance, an instance of $\beta * (\beta \text{ irr})$.

3.3 Context zipping

The intuition to think about these different cases is to consider that, for any σ' , we are looking for a way to construct a “witness” $\bar{\sigma}'$ such that $\tau[\bar{\sigma}'] = \sigma'$ from the hypothesis $\tau[\bar{\sigma}] \prec_v \sigma'$. When a type variable appears only once, its witness can be determined by inspecting the corresponding position in the type σ' . For example, in $\alpha * \beta \leq \text{bool} * \text{int}$, the mapping $\alpha \mapsto \text{bool}, \beta \mapsto \text{int}$ gives the witness pair bool, int .

However, when a variable appears twice, the two witnesses corresponding to the two occurrences may not coincide. (Consider for example $\beta * \beta \leq \text{bool} * \text{int}$.) If a variable β_i appears in several *invariant* occurrences, the witness of each occurrence is forced to be equal to the corresponding subterm of $\tau[\bar{\sigma}]$, that is σ_i , and therefore the various witnesses are themselves equal, hence compatible. On the contrary, for two covariant occurrences (as in the $\beta * \beta$ case), it is possible to pick a σ' such that the two witnesses are incompatible—and similarly for one covariant and one invariant occurrence. Finally, an irrelevant occurrence will never break closure properties, as all witnesses (forced by another occurrence) are compatible.

To express these merging properties, we define a *zip* operation $v_1 \hat{\wedge} v_2$, that formally expresses which combinations of variances are possible for several occurrences of the same variable; it is a partial operation (for example, it is not defined in the covariant-covariant case, which breaks the closure properties) with the following table:

| | | | | | |
|--------------------|-----|-----|-----|-----------|-----|
| $v \hat{\wedge} w$ | $=$ | $+$ | $-$ | \bowtie | w |
| $=$ | $=$ | | | $=$ | |
| $+$ | | | | $+$ | |
| $-$ | | | | $-$ | |
| \bowtie | $=$ | $+$ | $-$ | \bowtie | |
| v | | | | | |

3.4 Syntactic decomposability

Equipped with the zipping operation, we introduce a judgment $\Gamma \vdash \tau : v \Rightarrow v'$ to express decomposability, syntactically, defined by the inference rules on Figure 5. We also define its semantic interpretation $\llbracket \Gamma \vdash \tau : v \Rightarrow v' \rrbracket$. The judgment and its interpretation were co-designed, so keeping the interpretation in mind is the best way to understand the subtleties of the inference rules. We use zipping, which requires correct variances, to merge sub-derivations into larger ones, so, in addition to decomposability, the interpretation also ensures that v is a correct

$$\begin{array}{c}
\text{SC-TRIV} \\
\frac{v \geq v' \quad \Gamma \vdash \tau : v}{\Gamma \vdash \tau : v \Rightarrow v'} \\
\\
\text{SC-VAR} \\
\frac{w\alpha \in \Gamma \quad w = v}{\Gamma \vdash \alpha : v \Rightarrow v'} \\
\\
\text{SC-CONSTR} \\
\frac{\Gamma \vdash \mathbf{type} \overline{w\alpha} \mathbf{t} : v\text{-closed} \quad \Gamma = \hat{\lambda}_i \Gamma_i \quad \forall i, \Gamma_i \vdash \sigma_i : v.w_i \Rightarrow v'.w_i}{\Gamma \vdash \overline{\sigma} \mathbf{t} : v \Rightarrow v'}
\end{array}$$

Fig. 5. Syntactic decomposability

variance for τ under Γ . This subtlety is why we have two different properties for decomposability, $\Gamma \Vdash \tau : v \rightsquigarrow v'$ and $\llbracket \Gamma \vdash \tau : v \Rightarrow v' \rrbracket$.

Definition 6 (Interpretation of syntactic decomposability).

We write $\llbracket \Gamma \vdash \tau : v \Rightarrow v' \rrbracket$ for the conjunction of properties $\llbracket \Gamma \vdash \tau : v \rrbracket$ and $\Gamma \Vdash \tau : v \rightsquigarrow v'$.

To understand the inference rules, the first thing to notice is that the present rules are not completely syntax-directed: we first check whether $v \geq v'$ holds, and if not, we apply syntax-directed inference rules; existence of derivations is still easily decidable. If $v \geq v'$ holds, satisfying $\Gamma \Vdash \tau : v \rightsquigarrow v'$ (Definition 5) is trivial: $\tau[\overline{\sigma}] \prec_v \tau'$ implies $\tau[\overline{\sigma}] \prec_{v'} \tau'$, so taking $\overline{\sigma}$ for $\overline{\sigma}'$ is always a correct witness, which is represented by Rule SC-TRIV. The other rules then follow the same structure as the variance-checking judgment.

Rule SC-VAR is very similar to VC-VAR, except that the condition $w \geq v$ is replaced by a stronger equality $w = v$. This difference comes from the fact that the semantic condition for closure checking (Definition 2) includes both a variance check, which is monotonic in the context (Lemma 2) and the decomposability property, which is anti-monotonic (Lemma 4), so the present judgment must be invariant with respect to the context.

The most interesting rule is SC-CONSTR. It checks first that the head type constructor is v -closed (according to Definition 1); then, it checks that each subtype is decomposable from v to v' , *with compatible witnesses*, that is, in an environment family $(\Gamma_i)_{i \in I}$ that can be zipped into a unique environment Γ .

Lemma 5 (Soundness of syntactic decomposability).

If the judgment $\Gamma \vdash \tau : v \Rightarrow v'$ holds, then $\llbracket \Gamma \vdash \tau : v \Rightarrow v' \rrbracket$ holds.

Completeness in the general case is however much more difficult and we only prove it when the right-hand side variance v' is ($=$). In other words, we take back the generality that we have introduced in §3.1 when defining decomposability.

Lemma 6 (Completeness of syntactic decomposability).

If $\llbracket \Gamma \vdash \tau : v \Rightarrow v' \rrbracket$ holds for $v' \in \{=, \bowtie\}$, then $\Gamma \vdash \tau : v \Rightarrow v'$ is provable.

Lemma 6 is an essential piece to finally turn the semantic criterion REQ into a purely syntactic form.

Theorem 2 (Algorithmic criterion). *Given a variance annotation $(v_i \alpha_i)_{i \in I}$ and a constructor declaration of type $(\exists \bar{\beta} [\bigwedge_{i \in I} \alpha_i = T_i[\bar{\beta}]] . \tau[\bar{\beta}])$, the soundness criterion REQ for this constructor is equivalent to*

$$\exists \Gamma, (T_i)_{i \in I}, \quad \Gamma \vdash \tau : (+) \quad \wedge \quad \Gamma = \bigwedge_{i \in I} \Gamma_i \quad \wedge \quad \forall i \in I, \Gamma_i \vdash T_i : v_i \Rightarrow (=)$$

The three parts of this formula can be explained to a user, as soon as the underlying semantic phenomena (variable interference through zipping, and upward- and downward-closure) have been understood—there is no way to get around that. They are best read from right to left. The last part on the $(T_i)_{i \in I}$ is the decomposability requirement that failed in our example with $\langle m : \text{int} \rangle$: the type expressions equated with a covariant variable should be upward-closed, and those equated with a contravariant one downward-closed. The zipping part checks that the equations do not create interference through shared existential variables, as in **type** $(+\alpha, =\beta)$ **eq** = **Ref1 of** $\exists \gamma [\alpha = \gamma, \beta = \gamma]$. Finally, the variance check corresponds to the classic variance check on argument types of ADTs. One can verify that in presence of a simple ADT, this new criterion reduces to the simple syntactic criterion.

This presentation of the correctness criterion only relies on syntactic judgments. It is pragmatic in the sense that it suggests a simple and direct implementation, as a generalization of the check currently implemented in type system engines—which corresponds to the $\Gamma \vdash \tau : (+)$ part.

To compute the contexts Γ and $(T_i)_{i \in I}$ existentially quantified in this formula, one can use a variant of our syntactic judgments where the environment Γ is not an input, but an output of the judgment; in fact, one should return for each variable α the *set* of possible variances for this judgment to hold. For example, the query $(? \vdash \alpha * \beta \text{ ref} : +)$ should return $(\alpha \mapsto \{+, =\}; \beta \mapsto \{=\})$. Defining those algorithmic variants of the judgments is routine. The sets of variances corresponding to the decomposability of the $(T_i)_{i \in I}$ $(? \vdash T_i : v_i \Rightarrow (=))$ should be zipped together and intersected with the possible variances for τ , returned by $(? \vdash \tau : +)$. The algorithmic criterion is satisfied if and only if the intersection is not empty; this can be decided in a simple and efficient way.

4 Discussion

4.1 Upward and downward closure in a ML type system

In the type system we have used so far, all type constructors but p and q are both upward and downward-closed. This simple situation, however, does not hold in general: richer subtyping relations will have weaker invertibility properties. As soon as a bottom type \perp is introduced, for example, such that that for all type σ we have $\perp \leq \sigma$, downward-closure fails for all types – but \perp itself. For example, products are no longer downward-closed: $\Gamma \vdash \sigma * \tau \geq \perp$ does not implies that \perp is equal to some $\sigma' * \tau'$. Conversely, if one adds a top type \top , bigger than all other types, then most type are not upward-closed anymore.

In OCaml, there is no \perp or \top type⁸. However, object types and polymorphic variants have subtyping, so they are, in general, neither upward nor downward-closed. Finally, subtyping is also used in private type definitions, which were demonstrated in the example. Our closure-checking relation therefore degenerates into the following, quite unsatisfying, picture:

- no type is downward-closed because of the existence of private types;
- no object type but the empty object type is upward-closed;
- no arrow type is upward-closed because its left-hand-side would need to be downward-closed;
- datatypes are upward-closed if their components types are.

From a pragmatic point of view, the situation is not so bad; as our main practical motivation for finer variance checks is the relaxed value restriction, we care about upward-closure (covariance) more than downward-closure (contravariance). This criterion tells us that covariant parameters can be instantiated with covariant datatypes defined from sum and product types (but no arrow), which would satisfy a reasonable set of use cases.

4.2 A better control on upward and downward-closure

There is a subtle design question here. Decomposability is fundamentally a **negative** statement on the subtyping relation, guaranteeing that some types have no supertypes of a different structure. It is therefore not necessarily preserved by addition to the subtyping relation – our system, informally, is **non-monotone** in the subtyping relation.

This means that if we adopt the correctness criterion above, we must be careful in the future not to enrich the subtyping relation too much. Consider **private** types for example: one could imagine a symmetric concept of a type that would be strictly *above* a given type τ ; we will name those types **invisible** types (they can be constructed, but not observed). Invisible types and GADT covariance seem to be working against each other: if the designer adds one, adding the other later will be difficult.

A solution to this tension is to allow the user to *locally* guarantee negative properties about subtyping (what is *not* a subtype), at the cost of selectively abandoning the corresponding flexibility. Just as object-oriented languages have **final** classes that cannot be extended any more, we would like to be able to define some types as **downward-closed** (respectively **upward-closed**), that cannot later be made **private** (resp. **invisible**). Such declarations would be rejected if the defining type, for example an object type, already has subtypes (resp. supertypes), and would forbid further declarations of types below (resp. above) the defined type, effectively guaranteeing downward (resp. upward) closure.

⁸ A bottom type would be admissible, but a top type would be unsound in OCaml, as different types may have different runtime representations. Existential types, that may mix values of different types, are constructed explicitly through a boxing step.

Finally, upward or downward closure is a semantic aspect of a type that we must have the freedom to publish through an interface: abstract types could optionally be declared `upward-closed` or `downward-closed`.

4.3 Subtyping constraints and variance assignment

We will now revisit our example of strongly typed expressions in the introduction. A simple way to get such a type to be covariant would be, instead of proving delicate, non-monotonic upward-closure properties on the tuple type involved in the equation $\alpha = \beta * \gamma$, to *change* this definition so that the resulting type is obviously covariant:

```

type +α exp =
  | Val of ∃β[α ≥ β].β
  | Int of [α ≥ int].int
  | Thunk of ∃βγ[α ≥ γ].β exp * (β → γ)
  | Prod of ∃βγ[α ≥ β * γ].β exp * γ exp

```

We have turned each equality constraint $\alpha = T[\bar{\beta}]$ into a subtyping constraint $\alpha \geq T[\bar{\beta}]$. For a type α' such that $\alpha \leq \alpha'$, we get by transitivity that $\alpha' \geq T[\bar{\beta}]$. This means that $\alpha \text{ exp}$ trivially satisfies the correctness criterion REQ. Formally, instead of checking $\Gamma \vdash T_i : v_i \Rightarrow (=)$, we are now checking $\Gamma \vdash T_i : v_i \Rightarrow (+)$, which is significantly easier to satisfy: when v_i is itself $+$ we can directly apply the SC-TRIV rule. Note that this only works in the easy direction: while $\Gamma \vdash T_i : (+) \Rightarrow (+)$ is easy to check, $\Gamma \vdash T_i : (+) \Rightarrow (-)$ is just as hard as $\Gamma \vdash T_i : (+) \Rightarrow (=)$. In particular, an equality ($\sigma = \sigma'$) is already equivalent to a pair of inequalities ($\sigma \leq \sigma' \wedge \sigma \geq \sigma'$).

While this different datatype gives us a weaker subtyping assumption when pattern-matching, we are still able to write the classic function `eval` : $\alpha \text{ exp} \rightarrow \alpha$, because the constraints $\alpha \geq \tau$ are in the right direction to get an α as a result.

```

let rec eval : α exp → α = function
  | Val β (v : β) -> (v :> α)
  | Int (n : int) -> (n :> α)
  | Thunk β γ ((v : β exp), (f : β → γ)) ->
    (f (eval v) :> α)
  | Prod β γ ((b : β exp), (c : γ exp)) ->
    ((eval b, eval c) :> α)

```

This variation on GADTs, using subtyping instead of equality constraints, has been studied by Emir *et al* [EKRY06] in the context of the $\mathbb{C}\sharp$ programming language—it is also expressible in Scala. However, using subtyping constraints in GADTs has important practical drawbacks in a ML-like language. While typed object-oriented programming languages tend to use explicit polymorphism and implicit subtyping, ML uses implicit polymorphism and explicit subtyping (when present). Thus in ML, equality constraints can be implicitly used while subtyping constraints must be explicitly used: unification-based inference favors bidirectional equality over unidirectional subtyping. This makes GADT definitions

based on single subtyping constraints less convenient to use, because of the corresponding syntactic burden, and this is probably the reason why the notion of GADTs found in functional languages use only equality constraints. Subtyping constraints need also be explicit in the type declaration, forcing the user out of the convenient “generalized codomain type” syntax.

Finally, weakening equality constraints into a subtyping constraint in one direction is not always possible; sometimes the strictly weaker expressivity of the type forbids important uses. One must then use an equality constraint, and use our decomposability-based reasoning to justify the variance annotation. Consider the following example:

```

type +α tree =
  | Node of ∃β[α = β list]. (β tree) list
let append : α tree * α tree → α tree = function
  | Node β1 (l1 : β1 tree list), Node β2 (l2 : β2 tree list) ->
    Node (List.append l1 l2)

```

We know that the two arguments of `append` have the same type α `tree`. When matching on the `Node` constructors, we learn that α is equal to both β_1 `list` and β_2 `list`, from which we can deduce that β_1 is equal to β_2 by non-irrelevance of `list`. The concatenation of the lists `l1` and `l2` type-checks because this equality holds. If we used a type system without the decomposability criterion, we would need to turn the constructor constraint into $\exists\beta[\alpha \geq \beta$ `list`] to preserve covariance of α `tree`. We wouldn't necessarily have β_1 and β_2 equal anymore, so `(List.append l1 l2)`, hence the definition of `append` would not type-check. We would need decomposability-based reasoning to deduce, from $\alpha \geq \beta$ `list` and the fact that `list` is upward-closed, that in fact $\alpha = \beta'$ `list` for some β' .

This demonstrates that single subtyping constraints and our novel decomposability check on equality constraints are of incomparable expressivity: each setting handles programs that the other cannot type-check. From a theoretical standpoint, we think there is value in exploring the combination of both systems: using subtyping constraints rather than equalities, but also using decomposability to deduce stronger equalities when possible.

Note that while our soundness result directly transposes to a type-system with decomposability conditions on subtyping rather than equality constraints, our completeness result is special-cased on equality constraints. Completeness in the case of subtyping constraints is an open question.

Related Work

Simonet and Pottier [SP07] have studied GADTs in a general framework $\text{HMG}(X)$, inspired by $\text{HM}(X)$. They were interested in type inference using constraints, so considered GADTs with arbitrary constraints rather than type equalities, and considered the case of subtyping with applications to information flow security in mind. Their formulation of the checking problem for datatype declarations,

as a constraint-solving problem, is exactly our semantic criterion and is not amenable to a direct implementation. Correspondingly, they did not encounter any of the new notions of upward and downward-closure and variable interference (zipping) discussed in the present work. They define a dynamic semantics and prove that this semantic criterion implies subject reduction and progress. However, we cannot directly reuse their soundness result as they work in a setting where all constructors are upward- and downward-closed (their subtyping relation is atomic). We believe this is only an artifact of their presentation and their proof should be easily extensible to our setting.

Emir, Kennedy, Russo and Yu [EKRY06] studied the soundness of an object-oriented calculus with subtyping constraints on classes and methods. Previous work [KR05] had established the correspondence between equality constraints on methods in an object-oriented style and GADT constraints on type constructors in functional style. Through this surprisingly non-obvious correspondence, their system matches our presentation of GADTs with subtyping constraints and easier variance assignment, detailed in §4.3. They provide several usage examples and a full soundness proof using a classic syntactic argument. However, they do not consider the more delicate notions of decomposability, and their system therefore cannot handle some of the examples presented here.

Future Work

Experiments with v -closure of type constructors as a new semantic property In a language with non-atomic subtyping such as OCaml, we need to distinguish v -closed and non- v -closed type constructors. This is a new semantic property that, in particular, must be reflected through abstraction boundaries: we should be able to say about an abstract type that it is v -closed, or not say anything.

How inconvenient in practice is the need to expose those properties to have good variance for GADTs? Will the users be able to determine whether they want to enforce v -closure for a particular type they are defining?

Completeness of variance annotations with domain information The way we present GADTs using equality constraints instead of the codomain syntax is well-known to practitioners, under the form of a “factoring” transformation where an arbitrary GADT is expressed as a simple ADT, using the equality GADT $(\alpha, \beta) \text{ eq}$ as part of the constructor arguments to reify equality information.

This transformation does not work anymore with our current notion of GADTs in presence of subtyping. Indeed, all we can soundly say about the equality type $(\alpha, \beta) \text{ eq}$ is that it must be invariant in both its parameters; using $(\alpha, T_i[\bar{\beta}]) \text{ eq}$ as part of a constructor type would force the parameter α to be invariant.

We think it would be possible to re-enable factoring by eq by considering *domain information*, that is, information on constraints that must hold for the type to be inhabited. If we restricted the subtyping rule with conclusion $\bar{\sigma} \tau \leq \bar{\sigma}' \tau$ to only cases where $\bar{\sigma} \tau$ and $\bar{\sigma}' \tau$ are inhabited—with a separate rule to conclude subtyping in the non-inhabited case—we could have a finer variance check, as we

would only need to show that the criterion SEQ holds between two instances of the inhabited domain, and not any instance. If we stated that the domain of the type (α, β) eq is restricted by the constraint $\alpha = \beta$, we could soundly declare the variance $(\bowtie\alpha, \bowtie\beta)$ eq on this domain—which no longer prevents from factoring out GADTs by equality types.

Conclusion

Checking the variance of GADTs is surprisingly more difficult (and interesting) than we initially thought. We have studied a novel criterion of upward and downward closure of type expressions and proposed a corresponding syntactic judgment that is easily implementable. We presented a core formal framework to prove both its correctness and its completeness with respect to a natural semantic criterion.

This closure criterion exposes important tensions in the design of a subtyping relation, for which we previously knew of no convincing example in the context of ML-derived programming languages. We have suggested new language features to help alleviate these tensions, whose convenience and practicality is yet to be assessed by real-world usage.

Considering extensions of GADTs in a rich type system is useful in practice; it is also an interesting and demanding test of one’s type system design.

References

- Abe06. Andreas Abel. Polarized subtyping for sized types. *Mathematical Structures in Computer Science*, 2006. Special issue on subtyping, edited by Healfdene Goguen and Adriana Compagnoni.
- EKRY06. Burak Emir, Andrew Kennedy, Claudio Russo, and Dachuan Yu. Variance and generalized constraints for C# generics. In *Proceedings of the 20th European conference on Object-Oriented Programming, ECOOP’06*, 2006.
- Gar04. Jacques Garrigue. Relaxing the value restriction. In *In International Symposium on Functional and Logic Programming, Nara, LNCS 2998*, 2004.
- KR05. Andrew Kennedy and Claudio V. Russo. Generalized algebraic data types and object-oriented programming. In *Proceedings of the 20th annual ACM SIGPLAN conference on Object-oriented programming, systems, languages, and applications*, 2005. URL: <http://research.microsoft.com/pubs/64040/gadtoop.pdf>.
- Pfe01. Frank Pfenning. Intensionality, extensionality, and proof irrelevance in modal type theory. In *16th IEEE Symposium on Logic in Computer Science (LICS 2001), 16-19 June 2001, Boston University, USA, Proceedings*, 2001.
- SP07. Vincent Simonet and François Pottier. A constraint-based approach to guarded algebraic data types. *ACM Transactions on Programming Languages and Systems*, 29(1), January 2007. URL: <http://doi.acm.org/10.1145/1180475.1180476>.
- SR. Gabriel Scherer and Didier Rémy. GADTs meet subtyping. Long version, available electronically. URL: <http://gallium.inria.fr/~remy/gadts/>.