

GADT meet subtyping

Gabriel Scherer and Didier Rémy

INRIA, Rocquencourt*

Abstract. While generalized abstract datatypes (GADT) 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.

1 Motivation

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 whose programming idioms rely on subtyping, in particular object-oriented languages. Another 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, which may not classify mutable values. Therefore, it is important for extensions of type definitions, such as GADT, 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 + $\alpha$  expr =  
  | Val :  $\alpha \rightarrow \alpha$  expr  
  | Int : int  $\rightarrow$  int expr  
  | Thunk :  $\forall \beta. \beta$  expr * ( $\beta \rightarrow \alpha$ )  $\rightarrow$   $\alpha$  expr  
  | Prod :  $\forall \beta \gamma. \beta$  expr *  $\gamma$  expr  $\rightarrow$  ( $\beta * \gamma$ ) expr
```

Is it safe to say that `expr` is covariant in its type parameter? It turns out that, using the subtyping relation of the OCaml type system, the answer is “yes”.

* Part of this work has been done at IRILL.

But, surprisingly to us, in a type system with a top type \top , the answer would be “no”.

The aim of this article is to present a sound and complete criterion to check soundness of parameter variance annotations, for use in a type-checker. We also discuss the apparent fragility of this criterion with respect to changes to the subtyping relation (*e.g.* the presence or absence of a top type, private types, *etc.*), and a different, more robust way to combine GADT and subtyping.

Note Due to space restriction, the present article is only a short version from which many details have been omitted. All proofs of results presented in this version along with auxiliary results, as well as further discussions, can be found in the longer version available online [SR].

Examples

Let us first explain why it is reasonable to say that α `expr` 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 α `expr` into one of type α' `expr`. Here is some pseudo-code¹ for the coercion function:

```
let coerce :  $\alpha$  expr  $\rightarrow$   $\alpha'$  expr = function
| Val (v :  $\alpha$ ) -> Val (v :>  $\alpha'$ )
| Int n -> Int n
| Thunk  $\beta$  (b :  $\beta$  expr) (f :  $\beta \rightarrow \alpha$ ) ->
  Thunk  $\beta$  b (fun x -> (f x :>  $\alpha'$ ))
| Prod  $\beta$   $\gamma$  ((b, c) :  $\beta$  expr *  $\gamma$  expr) ->
  (* 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 :>  $\beta'$  expr), (c :>  $\gamma'$  expr))
```

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 $\gamma \leq \gamma'$. This allows to conclude by casting at types β' `expr` and γ' `expr`, recursively.

Similarly, in the `Int` case, we know that α must be an `int` and therefore an `int expr` is returned. This is because we know that, in OCaml, no type is above `int`: if `int` \leq τ , 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*

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 .

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 (+ $\alpha, =\beta$) eq`, from any $\sigma \leq \tau$ we could subtype $(\sigma, \sigma) \text{ eq}$ into $(\tau, \sigma) \text{ eq}$, allowing to cast any value of type τ back into one of type σ , which is unsound in general.

As a counter-example, the following declaration is incorrect: the type $\alpha \text{ t}$ cannot be declared covariant.

```
type + $\alpha \text{ t}$  =
  | K :  $\langle m : \text{int} \rangle \rightarrow \langle m : \text{int} \rangle \text{ t}$ 
  let v = (K (object method m = 1 end) :> < > t)
```

This declaration uses the OCaml object type $\langle m : \text{int} \rangle$, 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 $\langle \rangle$, so the alleged covariance of `t`, if accepted by the compiler, would allow us to cast a value of type $\langle m : \text{int} \rangle \text{ t}$ into one of type $\langle \rangle \text{ t}$. However, from such a value, we could wrongly deduce an equality witness $(\langle \rangle, \langle m : \text{int} \rangle) \text{ eq}$ that allows to cast any empty object of type $\langle \rangle$ into an object of type $\langle m : \text{int} \rangle$, but this is unsound, of course!

```
let get_eq :  $\alpha \text{ t} \rightarrow (\alpha, \langle m : \text{int} \rangle) \text{ eq}$  = function
  | K _ -> Refl      (* locally  $\alpha = \langle m : \text{int} \rangle$  *)
let wrong :  $\langle \rangle \rightarrow \langle m : \text{int} \rangle =$ 
  let eq :  $(\langle \rangle, \langle m : \text{int} \rangle) \text{ eq}$  = 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 t = τ internally` may describe its interface as `type t = private τ` . This is a compromise between a type abbreviation and an abstract type: it is possible to cast a value of type `t` into one of type τ , but not, conversely, to construct a value of type `t` from one of type τ . In other words, `t` is a strict subtype of τ : we have $t \leq \tau$ but not $t \geq \tau$. 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

³ This counterexample is due to Jeremy Yallop.

⁴ This counterexample is due to Jacques Garrigue.

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`).

In this article, we formally show that these notions of upward and downward-closure are the key to a sound variance check for GADT. We start from the formal development of Simonet and Pottier [SP07], which provides a general soundness proof for a language with subtyping and a very general notion of GADT expressing arbitrary constraints—rather than only type equalities. By specializing their correctness criterion, we can express it in terms of syntactic checks for closure and variance, that are simple to implement in a type-checker.

The problem of non-monotonicity

There is a problem with those upward or downward closure assumptions: while they hold in core ML, with strong inversion theorems, they are non-monotonic properties: they are not necessarily preserved by extensions of the subtyping lattice. For example, OCaml has a concept of *private types*: a type specified by `type t = private τ` is a new semi-abstract type smaller than τ ($t \leq \tau$ but $t \not\geq \tau$), that can be defined a posteriori for any type. Hence, no type is downward-closed *forever*. That is, for any type τ , a new, strictly smaller type may always be defined in the future.

This means that closure properties of the OCaml type system are relatively weak: no type is downward-closed⁵ (so instantiated GADT parameters cannot be contravariant), and arrow types are not upward-closed as their domain should be downward-closed. Only purely positive algebraic datatypes are upward-closed. The subset of GADT declarations that can be declared covariant today is small, yet, we think, large enough to capture a lot of useful examples, such as α `expr` above.

Giving back the freedom of subtyping

It is disturbing that the type system should rely on non-monotonic properties: 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

⁵ Except types that are only defined privately in a module and not exported: they exist in a “closed world” and we can check, for example, that they are never used in a `private` type definition.

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 **public** (respectively **visible**), that cannot later be made **private** (resp. **invisible**). Such declarations would be rejected if the defining type already has subtypes (*e.g.* an object type), and would forbid further declarations of types below (resp. above) the defined type, effectively guaranteeing downward (resp. upward) closure. 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 **public** or **visible**.

Another approach: subtyping constraints

Getting fine variance properties out of GADT is difficult because they correspond to type equalities which, to a first approximation, use their two operands both positively and negatively. One way to get an easy variance check is to encourage users to *change* their definitions into different ones that are easier to check. For example, consider the following redefinition of α **expr** (in a speculative extension of OCaml with subtyping constraints):

```

type + $\alpha$  expr =
  | Val :  $\forall \alpha. \alpha \rightarrow \alpha$  expr
  | Int :  $\forall \alpha[\alpha \geq \text{int}]. \text{int} \rightarrow \alpha$  expr
  | Thunk :  $\forall \beta. \beta$  expr * ( $\beta \rightarrow \alpha$ )  $\rightarrow \alpha$  expr
  | Prod :  $\forall \alpha \beta \gamma[\alpha \geq \beta * \gamma]. (\beta$  expr *  $\gamma$  expr)  $\rightarrow \alpha$  expr

```

It is now quite easy to check that this definition is covariant, since all type equalities $\alpha = T_i[\beta]$ have been replaced by inequalities $\alpha \geq T_i[\beta]$ which are preserved when replacing α by a subtype $\alpha' \geq \alpha$ —we explain this more formally in §4.2. This variation on GADT, using subtyping instead of equality constraints, has been studied by Emir *et al* [EKRY06] in the context of the C# programming language.

But isn't such a type definition less useful than the previous one, which had a stronger constraint? We will discuss this choice in more detail in §4.2.

Related work

Simonet and Pottier [SP07] have studied GADT in a general framework HMG(X), inspired by HM(X). They were interested in typing inference using constraints, so considered GADT with arbitrary constraints rather than type equalities, and considered the case of subtyping with applications to information flow security in mind. We instantiate their general framework, which allows us to reuse their dynamic semantics and syntactic proofs of soundness, and concentrate only on

the static semantics, proving that we meet the requirements they impose on the parametrized models.

Their soundness criterion is formulated in very general terms as a constraint entailment problem. In contrast, our specialized study of the case of equality and subtyping led to a refined, more syntactic, criterion. This provides a more practically implementable check for type definitions, and reveals the design issues surrounding v -closed constructors that were not apparent in their work.

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 correspondence, their system matches our presentation of GADT with subtyping constraints and easier variance assignment, mentioned in the introduction (§1) and detailed in §4.2. They do not encounter the more delicate notion of upward and downward closure.

Those two approaches (subtyping constraints with easy variance check, stronger equality constraints with more delicate variance check) are complementary and have different convenience trade-offs. In their system with explicit-constraint definitions and implicit subtyping, the subtyping-constraint solution is the most convenient, while our ML setting provides incentives to study the other solution.

2 A formal setting

We define a core language for Algebraic Datatypes (ADT) and, later, Generalized Algebraic Datatypes (GADT), that is an instance of the parametrized HMG(X) system of Simonet and Pottier [SP07]. We refine their framework by using variances to define subtyping, but rely on their formal description for most of the system, in particular the static and dynamic semantics. We ultimately rely on their type soundness proof, by rigorously showing (in the next section) that their requirements on datatype definitions for this proof to hold are met in our extension with variances.

2.1 Atomic subtyping

Our type system defines a subtyping relation between ground types, parametrized by a reflexive transitive relation between base constant types (`int`, `bool`, etc.). Ground types consist of a set of base types \mathbf{b} , function types $\tau_1 \rightarrow \tau_2$, product types $\tau_1 * \tau_2$, and a set of algebraic datatypes $\bar{\sigma} \mathbf{t}$. (We write $\bar{\sigma}$ for a sequence of types $(\sigma_i)_{i \in I}$.) We use prefix notation for datatype parameters, as is the usage in ML. Datatypes may be user-defined by toplevel declarations of the form:

```

type  $\bar{\sigma} \mathbf{t} =$ 
  |  $K_1$  of  $\tau^1[\bar{\alpha}]$ 
  | ...
  |  $K_n$  of  $\tau^n[\bar{\alpha}]$ 

```

$$\begin{array}{c}
\frac{}{\sigma \leq \sigma} \qquad \frac{\sigma_1 \leq \sigma_2 \quad \sigma_2 \leq \sigma_3}{\sigma_1 \leq \sigma_3} \qquad \frac{\mathbf{b} \leq \mathbf{c}}{\mathbf{b} \leq \mathbf{c}} \qquad \frac{\sigma \geq \sigma' \quad \tau \leq \tau'}{\sigma \rightarrow \tau \leq \sigma' \rightarrow \tau'} \\
\\
\frac{\sigma \leq \sigma' \quad \tau \leq \tau'}{\sigma * \tau \leq \sigma' * \tau'} \qquad \frac{\mathbf{type} \ \overline{v\alpha} \ \mathbf{t} \quad \forall i, \sigma_i \prec_{v_i} \sigma'_i}{\overline{\sigma} \ \mathbf{t} \leq \overline{\sigma'} \ \mathbf{t}}
\end{array}$$

Fig. 1. Subtyping relation

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

The sequence $\overline{v\alpha}$ is a binding list of type variables α_i along with their *variance annotation* v_i , which is a marker among 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), 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 $\overline{v\alpha}$. The judgment $\mathbf{type} \ \overline{v\alpha} \ \mathbf{t}$ simply means that the type constructor \mathbf{t} has been previously defined with the variance annotation $\overline{v\alpha}$. Notice that the rules for arrow and product types can be subsumed by the rule for datatypes, if one consider them as special datatypes (with a specific dynamic semantics) of variance $(-, +)$ and $(+, +)$, respectively. For this reason, the following definitions will not explicitly detail the cases for arrows and products.

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

On the restriction of atomic subtyping Our typing relation reproduces a simplification that is present in the formulation of Simonet and Pottier: it is *atomic* in the sense that two non-constant type constructors in the subtyping relation are always identical. We are confident their proof (and then our formal setting) can be extended to cover the non-atomic case, but we have left this extension to future work.

Richer type systems, for example if they have bottom or top types or, in the case of the OCaml type system, `private` types and object types, have a non-atomic subtyping relations. To be able to extend our work to such settings, we have carefully marked each use of atomic subtyping in the formal development

with an hypothesis of v -closure, defined below. In the case of atomic subtyping, all types are v -closed.

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 $\bar{\sigma}'$.*

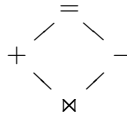
2.2 The algebra of variances

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 for each i , $\sigma_i \prec_{v_i} \sigma'_i$. 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 following table:

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

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 $\sigma \prec_v \tau$.⁶ This reflexive, partial order is described by the following lattice diagram:



That is, all variances are smaller than $=$ and bigger than \bowtie .

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 monotonous: if $v \leq v'$ then $w.v \leq w.v'$ (and $v.w \leq v'.w$).

⁶ 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** $w\alpha \mathbf{t}$: $\sigma \prec_w \sigma'$ 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 $\bowtie \leq +$.

We will frequently 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 for all i we have $\sigma_i \prec_{v_i} \sigma'_i$.

2.3 Variance assignment in ADTs

A counter-example To have a sound type system, some datatype declarations must be rejected. Assume (only for this example) that we have two base types `int` and `bool` such that `bool` \leq `int` and `int` $\not\leq$ `bool`. Consider the following type declaration:

```
type (+α, +β) t =
  | Fun of α → β
```

If it were accepted, we could build type the following program that deduces from the $(+\alpha)$ variance that `(bool, bool) t` \leq `(int, bool) t`; that is, we could turn the identity function of type `bool` \rightarrow `bool` into one of type `int` \rightarrow `bool` and then turns an integer into a boolean:

```
let three_as_bool : bool =
  match (Fun (fun x -> x) : (bool, bool) t :> (int, bool) t) with
  | Fun (danger : int → bool) -> danger 3
```

A requirement for type soundness We say that the type `type $\overline{v\alpha}$ t` defined by the constructors $(K_c \text{ of } \tau^c[\overline{\alpha}])_{c \in C}$ is *well-signed* if

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

The definition of $(+\alpha, +\beta) \mathbf{t}$ is not well-signed because we have $(\perp, \perp) \mathbf{t} \leq (\text{int}, \perp) \mathbf{t}$ according to the variance declaration, but we do not have the corresponding conclusion $(\text{int} \rightarrow \perp) \leq (\perp \rightarrow \perp)$.

This is a simplified version, specialized to simple algebraic datatypes, of the soundness criterion of Simonet and Pottier. They proved that this condition is *sufficient*⁷ for soundness: if all datatype definitions accepted by the type-checker are well-signed, then both subject reduction and progress hold—for their static and dynamic semantics, using the subtyping relation (\leq) we have defined.

A judgment for variance assignment When reformulating the well-signedness requirement of Simonet and Pottier for simple ADT, in our specific case where the subtyping relation is defined by variance, it becomes a simple check on the variance of type definitions. Our example above is unsound as it claims α covariant while it in fact appears in negative position in the definition.

⁷ It turns out that this condition is not *necessary* and can be slightly weakened: we discuss this in the extended version of this article.

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

Fig. 2. Variance assignment

We define a judgment to check the variance of a type expression. Given a context Γ of the form $\overline{w\alpha}$, 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 2.

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 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 semantics 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: $\forall \overline{\sigma}, \overline{\sigma}', \overline{\sigma} \prec_{\Gamma} \overline{\sigma}' \implies \tau[\overline{\sigma}] \prec_v \tau[\overline{\sigma}']$.

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 $\tau[\overline{\alpha}]$, variance v , and type sequences $\overline{\sigma}$ and $\overline{\sigma}'$, the subtyping relation $\tau[\overline{\sigma}] \prec_v \tau[\overline{\sigma}']$ holds if and only if the judgment $\Gamma \vdash \tau : v$ holds for some context Γ such that $\overline{\sigma} \prec_{\Gamma} \overline{\sigma}'$.

Furthermore, if $\tau[\overline{\sigma}] \prec_v \tau[\overline{\sigma}']$, then $\overline{\sigma} \prec_{\Delta} \overline{\sigma}'$ holds, where Δ is the more general context such that $\Delta \vdash \tau : v$ holds.

Checking variance of type definitions We have all the machinery in place to explain the checking of ADT variance declarations. The well-signedness criterion of Simonet and Pottier gives us a general semantic characterization of which definitions are correct: a definition $\mathbf{type} \ \bar{v}\bar{\alpha} \ \mathbf{t} = (\mathbf{K}_c \ \mathbf{of} \ (\tau^c[\bar{\alpha}]_{c \in C}))$ is correct if, for each constructor c , we have:

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

By inversion of subtyping, $\bar{\sigma} \ \mathbf{t} \leq \bar{\sigma}' \ \mathbf{t}$ implies $\sigma_i \prec_{v_i} \sigma'_i$ for all i . Therefore, it suffices to check that:

$$\forall \bar{\sigma}, \forall \bar{\sigma}', \quad (\forall i, \sigma_i \prec_{v_i} \sigma'_i) \implies \tau[\bar{\sigma}] \leq \tau[\bar{\sigma}']$$

This is exactly the semantic property corresponding to the judgment $\bar{v}\bar{\alpha} \vdash \tau : (+)$! That is, we have reduced soundness verification of an algebraic type definition to a mechanical syntactic check on the constructor argument type.

This syntactic criterion is very close to the one implemented in actual type checkers, which do not need to decide general subtyping judgments—or worse solve general subtyping constraints—to check variance of datatype parameters. Our aim is now to find a similar syntactic criterion for the soundness of variance annotations on guarded algebraic datatypes, rather than simple algebraic datatypes.

2.4 Variance annotations in GADT

A general description of GADT When used to build terms of type $\bar{\alpha} \ \mathbf{t}$, a constructor $\mathbf{K} \ \mathbf{of} \ \tau$ behaves like a function of type $\forall \bar{\alpha}. (\tau \rightarrow \bar{\alpha} \ \mathbf{t})$. Remark that the codomain is exactly $\bar{\alpha} \ \mathbf{t}$, the type \mathbf{t} instantiated with parametric variables. GADT arise by relaxing this restriction, allowing to specify 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:

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

Instead of being just $\alpha \ \mathbf{expr}$, the codomain is now $(\beta * \gamma) \ \mathbf{expr}$. 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 GADT support, such as Haskell and OCaml, and is reminiscent of the inductive datatype definitions of dependently typed languages.

However, for formal study of GADT, a different formulation based on equality constraints is preferred. We will use the following equivalent presentation, in the syntax of Simonet and Pottier. The idea is that instead of having $(\beta * \gamma) \ \mathbf{t}$ as codomain of the constructor `Prod`, we will force it to be $\alpha \ \mathbf{t}$ again, by adding an explicit type equality $\alpha = \beta * \gamma$.

```

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

```

In the rest of the paper, we extend our former core language with such guarded algebraic datatypes. This impacts the typing rules (which are precisely defined in Simonet and Pottier), but not the notion of subtyping, which is defined on (GADT) type constructors with variance `type $\overline{v\alpha}$ t` just as it previously was on simple datatypes. What needs to be changed, however, is the soundness criterion for checking the variance of type definitions.

The correctness criterion Simonet and Pottier [SP07] define a general framework HMG(X) to study type systems with GADT where the type equalities in bounded quantification are generalized to an arbitrary constraint language. They make few assumptions on the type system used, mostly that it has function types $\sigma \rightarrow \tau$, user-definable (guarded) algebraic datatypes $\overline{\alpha} \mathbf{t}$, and a subtyping relation $\sigma \leq \tau$ (which may be just equality, in languages without subtyping).

They use this general type system to give static semantics (typing rules) to a fixed untyped lambda-calculus equipped with datatype construction and pattern matching operations. They are able to prove a type soundness result under just some general assumptions on the particular subtyping relation (\leq). Here are the three requirements to get their soundness result:

1. Incomparability of distinct types: for all types $\tau_1, \tau_2, \overline{\sigma}, \overline{\sigma}'$ and distinct datatypes $\overline{\alpha} \mathbf{t}, \overline{\alpha}' \mathbf{t}'$, the types $(\tau_1 \rightarrow \tau_2), \tau_1 * \tau_2, \overline{\sigma} \mathbf{t}$ and $\overline{\sigma}' \mathbf{t}'$ must be pairwise incomparable (both $\not\leq$ and $\not\geq$) — this is where our restriction to an atomic subtyping relation, discussed in §2.1, comes from.
2. Decomposability of function and product types: if $\tau_1 \rightarrow \tau_2 \leq \sigma_1 \rightarrow \sigma_2$ (respectively $\tau_1 * \tau_2 \leq \sigma_1 * \sigma_2$), we must have $\tau_1 \geq \sigma_1$ (resp. $\tau_1 \leq \sigma_1$) and $\tau_2 \leq \sigma_2$.
3. Decomposability of datatypes⁸: for each datatype $\overline{\alpha} \mathbf{t}$ and all type vectors $\overline{\sigma}$ and $\overline{\sigma}'$ such that $\overline{\sigma} \mathbf{t} \leq \overline{\sigma}' \mathbf{t}$, we must have $(\exists\overline{\beta}[D[\overline{\sigma}]]\tau) \leq (\exists\overline{\beta}[D[\overline{\sigma}']]\tau)$ for each constructor `K of $\exists\overline{\beta}[D[\overline{\beta}, \overline{\alpha}]]. \tau[\overline{\beta}]$` .

Those three criteria are necessary for the soundness proof. We will now explain how variance of type parameters impact those requirements, that is, how to match a GADT implementation against a variance specification. With our definition of subtyping based on variance, and the assumption that the datatype `$\overline{v\alpha} \mathbf{t}$` we are defining indeed has variance `$\overline{v\alpha}$` , is the GADT decomposability requirement (item 3 above) satisfied by all its constructors? If so, then the datatype definition is sound and can be accepted. Otherwise, the datatype definition does not match the specified variance, and should be rejected by the type checker.

⁸ This is an extended version of the soundness requirement for algebraic datatypes: it is now formulated in terms guarded existential types $\exists\overline{\beta}[D]\tau$ rather than simple argument types τ .

3 Checking variances of GADT

For every type definition, we need to check that the decomposability requirement of Simonet and Pottier holds. Remark that it is expressed for each GADT constructor independently of the other constructors for the same type: we can check one constructor at a time.

Assume we check a fixed constructor K of argument type $\exists \bar{\beta}[D[\bar{\alpha}, \bar{\beta}]]. \tau[\bar{\beta}]$. Simonet and Pottier prove that their requirement is equivalent to the following formula, which is more convenient to manipulate:

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

The purpose of this section is to extract a practical criterion equivalent to this requirement. It should not be expressed as a general constraint satisfaction problem, but rather as a syntax-directed and decidable algorithm that can be used in a type-checker—without having to implement a full-blown constraint solver.

3.1 Expressing decomposability

If we specialize REQ-SP to the **Prod** constructor of the α **expr** example datatype, *i.e.* **Prod** of $\exists \beta \gamma [\alpha = \beta * \gamma] \beta$ **expr** * γ **expr**, we get:

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

We can substitute equalities and use the (assumed) covariance to simplify the subtyping constraint $\sigma \mathbf{expr} \leq \sigma' \mathbf{expr}$ 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. This 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 form

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

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:

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

Our aim is now to find a set of inference rules to check decomposability; we will later reconnect it to REQ-SP. 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 3 (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 property

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

holds.

We use the symbol \Vdash rather than \vdash to highlight the fact that this is just a logic formula, not the semantics criterion corresponding to an inductive judgment, nor 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-monotonous” 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.

In the following subsections, we study the subtleties of decomposability.

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—and any other head type constructor. In the general case, we require that this specific type constructor be upward-closed.

In the general case, 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`, but it is not equal to $(\beta * \beta)[\gamma']$ for any γ' . 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 instance⁹ of $\beta \text{ ref } * \beta \text{ ref}$.

⁹ 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”¹⁰ 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 3.

We sometimes need to merge sub-derivations into larger ones, so in addition to decomposability, the judgments simultaneously ensures that v is a correct

¹⁰ The idea of context merging and the term “zipping” are inspired by Montagu and Remy [MR09]

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

Fig. 3. Syntactic decomposability

variance for τ under Γ . Actually, in order to understand the details of this judgment, it is quite instructive to compare it with the variance-checking judgment $\Gamma \vdash \tau : v$ defined on Figure 2.

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 the semantics criterion 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 semantics 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 each subtype for decomposability from v to v' *with compatible witnesses*, that is, in an environment family Γ_i that can be zipped into a unique environment Γ .

In order to connect the syntactic and semantics versions of decomposability, we define the interpretation $\llbracket \Gamma \vdash \tau : v \Rightarrow v' \rrbracket$ of syntactic decomposability.

Definition 4 (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 $\Gamma \Vdash \tau : v \rightsquigarrow v'$.

Note that our interpretation $\Gamma \vdash \tau : v \Rightarrow v'$ does not coincide with our previous decomposability formula $\Gamma \Vdash \tau : v \rightsquigarrow v'$, because of the additional variance-checking hypothesis that makes it composable. The distinction between those two notions of decomposition is not useful to have a sound criterion, but is crucial to be complete with respect to the criterion of Simonet and Pottier, which imposes no variance checking condition.

Lemma 5 (Soundness of syntactic decomposability).

If the judgment $\Gamma \vdash \tau : v \Rightarrow v'$ holds, then $\llbracket \Gamma \vdash \tau : v \Rightarrow v' \rrbracket$ is 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 correctness criterion REQ-SP of Simonet and Pottier into a purely syntactic criterion.

Theorem 2 (Algorithmic criterion). *The REQ-SP criterion is equivalent to*

$$\exists \Gamma, (T_i)_i, \quad \Gamma \vdash \tau : (+) \quad \wedge \quad \Gamma = \bigwedge_i T_i \quad \wedge \quad \forall i, \Gamma_i \vdash T_i : v_i \Rightarrow (=)$$

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 are only concerned with 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, and we have not done it hereby lack of space. 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 Closed-world vs. open-world subtyping

4.1 A better control on upward and downward-closure

As explained in the introduction, the problem with the upward and downward closure properties is that they are not monotonic: enriching the subtyping lattice of our type system does not preserve them. While the core language has a nice variance check for GADT, adding private types in particular destroys the downward-closure property of the whole type system.

Our proposed solution to this tension is to give the user the choice to locally strengthen negative knowledge about the subtyping relation by abandoning some flexibility. Just as object-oriented languages have a concept of **final** classes that cannot be extended, we would like to allow to define **downward-closed** datatypes, whose private counterparts cannot be declared, and **upward-closed** datatypes that cannot be made **invisible**: defining **type t = private τ** would be rejected by the type-checker if τ was itself declared **downward-closed**.

4.2 Subtyping constraints and variance assignment

Consider our introductory example α **expr** of strongly typed expressions (§1). 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 + $\alpha$  expr =
  | Val of  $\exists\beta[\alpha \geq \beta].\beta$ 
  | Int of  $[\alpha \geq \text{int}].\text{int}$ 
  | Thunk of  $\exists\beta\gamma[\alpha \geq \gamma].\beta$  expr *  $(\beta \rightarrow \gamma)$ 
  | Prod of  $\exists\beta\gamma[\alpha \geq \beta * \gamma].\beta$  expr *  $\gamma$  expr

```

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 α **expr** trivially satisfies the correctness criterion from Simonet and Pottier. 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.

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

```

let rec eval :  $\alpha$  expr  $\rightarrow \alpha$  = function
  | Val  $\beta$  (v :  $\beta$ ) -> (v :>  $\alpha$ )
  | Int (n : int) -> (n :>  $\alpha$ )
  | Thunk  $\beta\gamma$  ((v :  $\beta$  expr), (f :  $\beta \rightarrow \gamma$ )) ->
    (f (eval v) :>  $\alpha$ )
  | Prod  $\beta\gamma$  ((b :  $\beta$  expr), (c :  $\gamma$  expr)) ->
    ((eval b, eval c) :>  $\alpha$ )

```

However, allowing subtyping constraints in GADT has some disadvantages. If the language requires subtyping casts to be explicit, this would make pattern matching of GADT syntactically heavier than with current GADT where equalities constraints are used implicitly. Subtyping constraints need also be explicit in the type declaration, forcing the user out of the convenient “generalized codomain type” syntax.

From a theoretical standpoint, we think there is value in exploring both directions: experimenting with GADT using subtyping constraints, and with fine-grained closure properties for equality constraints. Both designs allow to reason in an open world setting, by being resilient to extensions of the subtyping relation. Whether it is possible to expose those features to the expert language user (*e.g.* library designers) without forcing all users to pay the complexity burden remains to be seen.

5 Future Work

Extension of the formal exposition to non-atomic subtyping As remarked in §2.1 during the definition of our formal subtyping relation, the soundness proof of Simonet and Pottier is restricted to atomic subtyping. We conjecture that their work can be extended to non-atomic subtyping, and furthermore that our results would extend seamlessly in this setting, thanks to our explicit use of the v -closure hypothesis.

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 GADT? Will the users be able to determine whether they want to enforce v -closure for a particular type they are defining?

Experiments with subtyping constraints in GADT In §4.2, we have presented a different way to define GADT with weaker constraints (simple subtyping instead of equality) and stronger variance properties. It is interesting to note that, for the few GADT examples we have considered, using subtyping constraints rather than equality constraints was sufficient for the desired applications of the GADT.

However, there are cases where the strong equality relying on fine-grained closure properties is required. We need to consider more examples of both cases to evaluate the expressiveness trade-off in, for example, deciding to add only one of these solutions to an existing type system.

Completeness of variance annotations with domain information For simple algebraic datatypes, variance annotations are “enough” to say anything we want to say about the variance of datatypes. Essentially, all admissible variance relations between datatypes can be described by considering the pairwise variance of their parameters, separately.

This does not work anymore with GADT. For example, with only this notion of variance, all we can soundly say about the equality type (α, β) **eq** is that it must be invariant in both its parameters. In particular, the well-known trick of “factoring out” GADT by using the **eq** type in place of equality constraint does not preserve variances.

We think it would be possible to regain some “completeness”, and in particular 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 $\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}$ to only cases where $\bar{\sigma} \mathbf{t}$ and $\bar{\sigma}' \mathbf{t}$ 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 of Simonet and Pottier 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 $(\varkappa\alpha, \varkappa\beta)$ `eq` on this domain—which no longer prevents from factoring out GADT by equality types.

Conclusion

Checking the variance of GADT 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 the more general criterion of Simonet and Pottier.

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 extension of GADT 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>.
- MR09. Benoît Montagu and Didier Rémy. Modeling abstract types in modules with open existential types. In *ACM Symposium on Principles of Programming Languages (POPL)*, pages 63–74, January 2009. URL: <http://gallium.inria.fr/~remy/modules/Montagu-Remy@popl09:fzip.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.
- SR. Gabriel Scherer and Didier Rémy. Gadt meet subtyping. Long version, available electronically. URL: <http://gallium.inria.fr/~remy/gadts/>.