



## GADT meet Subtyping

Gabriel Scherer, Didier Rémy

► **To cite this version:**

Gabriel Scherer, Didier Rémy. GADT meet Subtyping. [Research Report] RR-8114, INRIA. 2012, pp.33. <hal-00744292>

**HAL Id: hal-00744292**

**<https://hal.inria.fr/hal-00744292>**

Submitted on 22 Oct 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# GADT meet subtyping

Gabriel Scherer, Didier Rémy

**RESEARCH  
REPORT**

**N° 8114**

October 2012

Project-Team Gallium

ISRN INRIA/RR--8114--FR+ENG

ISSN 0249-6399





## GADT meet subtyping

Gabriel Scherer, Didier Rémy

Project-Team Gallium

Research Report n° 8114 — October 2012 — 33 pages

**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.

**Key-words:** subtyping, datatypes, variance

---

Part of this work has been done at IRILL

**RESEARCH CENTRE  
PARIS – ROCQUENCOURT**

Domaine de Voluceau, - Rocquencourt  
B.P. 105 - 78153 Le Chesnay Cedex

## GADT avec sous-typage

**Résumé :** Les types algébriques généralisés (*Generalized Algebraic Datatypes*, GADT) sont maintenant bien compris, mais leur ajout à un langage équipé de sous-typage nous réservait quelques surprises. Qu'est-ce qu'être covariant pour un paramètre de GADT ? La réponse s'avère difficile. Elle met en jeu des propriétés fines de la relation de sous-typage qui soulèvent d'intéressantes problématiques de conception de langage. Nous permettons des annotations de variance dans les définitions de GADT, étudions leur correction, et présentons un algorithme correct et complet pour les vérifier. Notre travail peut s'appliquer à un langage complet inspiré de ML et avec sous-typage explicite, tel que OCaml, ou même à des langages avec des contraintes générales de sous-typage.

**Mots-clés :** sous-typage, types de données, variance

## 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  $\alpha$  `list` of immutable lists is expected to be *covariant*: we wish  $\sigma$  `list`  $\leq$   $\sigma'$  `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. This relaxation uses an intuitive subtyping argument: all occurrences of such type variables can be specialized to  $\perp$ , and any-time later, all covariant occurrences of the same variable (which are now  $\perp$ ) can be simultaneously replaced by the same arbitrary type  $\tau$ , which is always a supertype of  $\perp$ . This relaxation of the value-restriction is implemented in OCaml, where it is surprisingly useful. 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”. 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.

## Examples

Let us first explain why it is reasonable to say that  $\alpha$  `expr` is covariant. Informally, if we are able to coerce a value of type  $\alpha$  into one of type  $\alpha'$  (we write  $(v :> \alpha')$  to explicitly cast a value  $v$  of type  $\alpha$  to a value of type  $\alpha'$ ), then we are also able to transform a value of type  $\alpha$  `expr` into one of type  $\alpha'$  `expr`. Here is some pseudo-code<sup>1</sup> 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  $\gamma \leq \gamma' *$  *)
```

<sup>1</sup>The variables  $\beta'$  and  $\gamma'$  of the `Prod` case are never really defined, only justified at the meta-level, making this code only an informal sketch.

```
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  $\alpha$  must be equal to some type of the form  $\beta * \gamma$ . So from  $\alpha \leq \alpha'$  we know that  $\beta * \gamma \leq \alpha'$ . Therefore,  $\alpha'$  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  $\beta'$  `expr` and  $\gamma'$  `expr`, recursively.

Similarly, in the `Int` case, we know that  $\alpha$  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`  $\leq \tau$ , then  $\tau$  must be `int`.

What we use in both cases is reasoning of the form<sup>2</sup>: “if  $T[\bar{\beta}] \leq \alpha'$ , then I know that  $\alpha'$  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  $\alpha'$  is of the form  $T[\bar{\beta}']$ .

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

```
type ( $\alpha$ ,  $\beta$ ) eq =
  | Refl :  $\forall \gamma. (\gamma, \gamma)$  eq

let cast (eqab : ( $\alpha$ ,  $\beta$ ) eq) :  $\alpha \rightarrow \beta$  =
  match eqab with
  | Refl -> (fun x -> x)
```

Notice that it would be unsound<sup>3</sup> 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  $(\sigma, \sigma)$  `eq` into  $(\tau, \sigma)$  `eq`, allowing to cast any value of type  $\tau$  back into one of type  $\sigma$ , which is unsound in general.

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

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

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 `t`, if accepted by the compiler, would allow us to cast a value of type `< m : int > t` into one of type `< > t`. However, from such a value, we could wrongly deduce 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 :  $\alpha$  t  $\rightarrow$  ( $\alpha$ , < 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*<sup>4</sup>: a module using a type `type t =  $\tau$  internally` may describe its interface as `type t = private  $\tau$` . This is a compromise

<sup>2</sup>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  $T$ .

<sup>3</sup>This counterexample is due to Jeremy Yallop.

<sup>4</sup>This counterexample is due to Jacques Garrigue.

between a type abbreviation and an abstract type: it is possible to cast a value of type  $\mathbf{t}$  into one of type  $\tau$ , but not, conversely, to construct a value of type  $\mathbf{t}$  from one of type  $\tau$ . In other words,  $\mathbf{t}$  is a strict subtype of  $\tau$ : we have  $\mathbf{t} \leq \tau$  but not  $\mathbf{t} \geq \tau$ . 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`.

Unsound GADT covariance declarations would defeat the purpose of such private types: as soon as the user gets one element of the private type, she could forge values of this type, as illustrated by the code below.

```

module M = struct
  type file_descr = int
  let stdin = 0
  let open = ...
end : sig
  type file_descr = private int
  val stdin : file_descr
  val open : string -> (file_descr, error) sum
end

type + $\alpha$  t =
  | K : priv -> M.file_descr t

let get_eq :  $\alpha$  t -> ( $\alpha$ , M.file_descr) eq = function
  | K _ -> Refl

let forge : int -> M.file_descr =
  fun (x : int) -> cast (get_eq p) M.stdin

```

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 `M.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  $\tau$`  is a new semi-abstract type smaller than  $\tau$  ( $\mathbf{t} \leq \tau$  but  $\mathbf{t} \not\geq \tau$ ), that can be defined a posteriori for any type. Hence, no type is downward-closed *forever*. That is, for any type  $\tau$ , 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-closed<sup>5</sup> (so instantiated GADT parameters cannot be

<sup>5</sup>Except types that are only defined privately in a module and not exported: they exist in



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  $\alpha$  `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  $\tau$ ; 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 `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  $\alpha$  `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  $\alpha$  by a subtype  $\alpha' \geq \alpha$ —we explain this more formally in §4.3. This variation on GADT, using subtyping instead of equality constraints, has been studied by Emir *et al* [EKRY06] in the context of the  $\mathbb{C}\sharp$  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 §4.3.

---

a “closed world” and we can check, for example, that they are never used in a `private` type definition.

## On the importance of variance annotations

Being able to specify the variance of a parametrized datatype is important at abstraction boundaries: one may wish to define a program component relying on an *abstract* type, but still make certain subtyping assumptions on this type. Variance assignments provide a framework to specify such a semantic interface with respect to subtyping. When this abstract type dependency is provided by an encapsulated implementation, the system must check that the provided implementation indeed matches the claimed variance properties.

Assume the user specifies an abstract type

```
module type S = sig
  type (+α) collection
  val empty : unit -> α collection
  val app : α collection -> α collection -> α collection
end
```

and then implements it with linked lists

```
module C : S = struct
  type +α collection =
    | Nil of unit
    | Cons of α * α collection
  let empty () = Nil ()
end
```

The type-checker will accept this implementation, as it has the specified variance.

On the contrary,

```
type +α collection = (α list) ref
let empty () = ref []
```

would be rejected, as `ref` is invariant. In the following definition:

```
let nil = C.empty ()
```

the right hand-side is not a value, and is therefore not generalized in presence of the value restriction; we get a monomorphic type,  $?α \mathbf{t}$ , where  $?α$  is a yet-undetermined type variable. The relaxed value restriction [Gar04] indicates that it is sound to generalize  $?α$ , as it only appears in covariant positions. Informally, one may unify  $?α$  with  $\perp$ , add an innocuous quantification over  $α$ , and then generalize  $\forall α. \perp \mathbf{t}$  into  $\forall α. α \mathbf{t}$  by covariance—assuming a lifting of subtyping to polymorphic type schemes.

The definition of `nil` will therefore get generalized in presence of the relaxed value restriction, which would not be the case if the interface `S` had specified an invariant type.

## Related work

When we encountered the question of checking variance annotations on GADT, we expected to find it already discussed in the literature. The work of Simonet and Pottier [SP07] is the closest we could find. It was done in the context of finding good specification for *type inference* of code using GADT, and in this context it is natural to embed some form of constraint solving in the type inference problem. From there, Simonet and Pottier generalized to a rich notion of GADT defined over arbitrary constraints, in presence of a subtyping relation, justified in their setting by potential applications to information flow checking.

They do not describe a particular type system, but a parametrized framework  $HMG(X)$ , in the spirit of the inference framework  $HM(X)$ . In this setting, they prove

a general soundness result, applicable to all type systems which satisfy their model requirements. We directly reuse this soundness result, by checking that we respect these requirements and proving that their condition for soundness is met. This allows us to concentrate purely on the static semantics, without having to define our own dynamic semantics to formulate subject reduction and progress results.

Their soundness requirement is formulated in terms of a general constraint entailment problem involving arbitrary constraints. Specializing this to our setting is simple, but expressing it in a form that is amenable to mechanical verification is surprisingly harder—this is the main result of this paper. Furthermore, at their level of generality, the design issues related to subtyping of GADT, in particular the notion of upward and downward-closed type constructors, were not apparent. Our article is therefore not only a specialized, more practical instance of their framework, but also raises new design issues.

The other major related work, by Emir, Kennedy, Russo and Yu [EKRY06], studies the soundness of having subtyping constraints on classes and methods of an object-oriented type system with generics (parametric polymorphism). Previous work [KR05] had already established the relation between the GADT style of having type equality constraints on data constructors and the desirable object-oriented feature of having type equality constraints on object methods. This work extends it to general subtyping constraints and develops a syntactic soundness proof in the context of a core type system for an object-oriented languages with generics.

The general duality between the “sums of data” prominent in functional programming and “record of operations” omnipresent in object-oriented programming is now well-understood. Yet, it is surprisingly difficult to reason on the correspondence between GADT and generalized method constraints; an application that is usually considered to require GADT in a functional style (for example a strongly-typed `eval`  $\alpha$  `expr` datatype and its associated `eval` function) is simply expressed in idiomatic object-oriented style without specific constraints<sup>6</sup>, while the simple `flatten`  $: \forall \alpha, \alpha$  `list list`  $\rightarrow$   $\alpha$  `list` requires an equality or subtyping constraint when expressed in object-oriented style.

These important differences of style and applications make it difficult to compare our present work with this one. Our understanding of this system is that a subtyping constraint of the form  $X \leq Y$  is considered to be a negative occurrence of  $X$ , and a positive occurrence of  $Y$ ; this means that equality constraints (which are conjunctions of a ( $\leq$ ) constraint and a ( $\geq$ ) constraints) always impose invariance on their arguments. Checking correctness of constraints with this notion of variance is simpler than with our upward and downward-closure criterion, but also not as expressive. It corresponds, in our work, to the idea of GADT with subtyping constraint mentioned in the introduction and that we detail in §4.3.

The design trade-off in this related work is different from our setting; the reason why we do not focus on this solution is that it requires explicit annotations at the GADT definition site, and more user annotations in pattern matching in our system where subtyping is explicitly annotated, while convertibility is implicitly inferred by unification. On the contrary, in an OOP system with explicit constraints and implicit subtyping, this solution has the advantage of user convenience.

We can therefore see our present work as a different choice in the design space: we want to allow a richer notion of variance assignment for type *equalities*, at the cost a higher complexity for those checks. Note that the two directions are complementary and are both expressed in our formal framework.

<sup>6</sup>There is a relation between this way of writing a strongly typed `eval` function and the “finally tagless” approach [Kis] that is known to require only simple ML types.

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

Figure 1: Subtyping relation

## 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  $\overline{\sigma} \mathbf{t}$ . (We write  $\overline{\sigma}$  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:

$$\begin{array}{l}
\text{type } \overline{v\alpha} \mathbf{t} = \\
\quad | K_1 \text{ of } \tau^1[\overline{\alpha}] \\
\quad | \dots \\
\quad | K_n \text{ of } \tau^n[\overline{\alpha}]
\end{array}$$

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  $\alpha_i$  along with their *variance annotation*  $v_i$ , which is a marker among the set  $\{+, -, =, \bowtie\}$ . We may associate a relation 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  $\sigma$  and  $\tau$ .

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  $\text{type } \overline{v\alpha} \mathbf{t}$  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.

Finally, it is routine to show that the rules for reflexivity and transitivity are admissible, by pushing them up in the derivation until the base cases  $\mathbf{b} \leq \mathbf{c}$ , where they can be removed as  $(\leq)$  is assumed to be reflexive and transitive. Removing reflexivity and transitivity provides us with an equivalent syntax-directed judgment having powerful inversion principles: 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  $i$ ,  $\sigma_i \prec_{v_i} \sigma'_i$ .

We insist that our equality relation  $(=)$  is here a derived concept, defined from the subtyping relation  $(\leq)$  as the “equiconvertibility” relation  $(\leq \cap \geq)$ ; in particular, it is not defined as the usual syntactic equality. If we have both  $b_1 \leq b_2$  and  $b_2 \leq b_1$  in our relation on base types, for two distinct base types  $b_1$  and  $b_2$ , we have  $b_1 = b_2$  as types, even though they are syntactically distinct. This choice is inspired by the previous work of Simonet and Pottier.

**On the restriction of atomic subtyping** The subtyping system demonstrated above is called “atomic”. If two head constructors are in the subtyping relation, they are either identical or constant (no parameters). Structure-changing subtyping occurs only at the leaves of the subtyping derivations.

While this simplifies the meta-theoretic study of the subtyping relation, this is too simplifying for real-world type systems that use non-atomic subtyping relations. In our examples using the OCaml type system, `private` type were a source of non-atomic subtyping: if you define `type  $\alpha$   $\mathbf{t}_2 = \text{private } \alpha \mathbf{t}_1$` , the head constructors  $\mathbf{t}_1$  and  $\mathbf{t}_2$  are distinct yet in a subtyping relation. If we want to apply our formal results to the design of such languages, we must be careful to isolate any assumption on this atomic nature of our core formal calculus.

The aspect of non-atomic subtype relations we are interested in is the notion of  $v$ -closed constructor. We have used this notion informally in the first section (in OCaml, product types are  $+$ -closed); we now defined it formally.

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

In our core calculus, all type constructors are  $v$ -closed for any  $v \neq \bowtie$ , but we will still mark this hypothesis explicitly when it appears in typing judgments; this let the formal results be adapted more easily to a non-atomic type system.

It would have been even more convincing to start from a non-atomic subtyping relation. However, the formal system of Simonet and Pottier, whose soundness proof we ultimately reuse, restricts subtyping relations between (G)ADT type to atomic subtyping. 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.

## 2.2 The algebra of variances

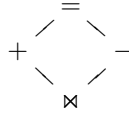
If we know that  $\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}$ , that is  $\bar{\sigma} \mathbf{t} \prec_+ \bar{\sigma}' \mathbf{t}$ , and the constructor  $\mathbf{t}$  has variable  $\bar{v}\bar{\alpha}$ , an inversion principle tells us that for each  $i$ ,  $\sigma_i \prec_{v_i} \sigma'_i$ . 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  $\sigma$  and  $\tau$ ,  $\sigma \prec_w \tau$  implies  $\sigma \prec_v \tau$ .<sup>7</sup> 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  $v.w \leq v'.w$ ).

We will frequently manipulate vectors  $\overline{v\alpha}$ , of variable associated with variances, which correspond to the “context”  $\Gamma$  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  $(\prec_\Gamma)$  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 (+ $\alpha$ , + $\beta$ ) t =
  | Fun of  $\alpha \rightarrow \beta$ 

```

If it were accepted, we could build 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

```

<sup>7</sup>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`  $v\alpha$  `t`, it is safe to assume `type`  $w\alpha$  `t`:  $\sigma \prec_w \sigma'$  implies  $\sigma \prec_v \sigma'$ , which implies  $\sigma t \leq \sigma' t$ . 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  $\bowtie \leq +$ .

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

Figure 2: Variance assignment

**A requirement for type soundness** We say that the type  $\mathbf{type} \overline{w\alpha} \mathbf{t}$  defined by the constructors  $(\mathbb{K}_c \text{ of } \tau^c[\overline{\alpha}])_{c \in C}$  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  $(\perp, \perp) \mathbf{t} \leq (\mathbf{int}, \perp) \mathbf{t}$  according to the variance declaration, but we do not have the corresponding conclusion  $(\mathbf{int} \rightarrow \perp) \leq (\perp \rightarrow \perp)$ .

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*<sup>8</sup> 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  $\alpha$  covariant while it in fact appears in negative position in the definition.

In the context of higher-order subtyping [Abe06], where type abstractions are first-class and annotated with a variance  $(\lambda v\alpha.\tau)$ , it is natural to present this check as a kind checking of the form  $\Gamma \vdash \tau : \kappa$ , where  $\Gamma$  is a context  $\overline{w\alpha}$  of type variables *associated with variances*. For example, if  $+\alpha \vdash \tau : \star$  is provable, it is sound to consider  $\alpha$  covariant in  $\tau$ . In the context of a simple first-order monomorphic type calculus, this amounts to a *monotonicity check* on the type  $\tau$  as defined by [EKRY06]. Both approaches use judgments of a peculiar form where the *context* changes when going under a type constructor: to check  $\Gamma \vdash \sigma \rightarrow \tau$ , one checks  $\Gamma \vdash \tau$  but  $(\Gamma/-) \vdash \sigma$ , where  $\Gamma/-$  reverses all the variances in the context  $\Gamma$  (turns  $(-)$  into  $(+)$  and conversely). Abel gives an elegant presentation of this inversion  $/$  as an algebraic operation on variances, a quasi-inverse such that  $u/v \leq w$  if and only if  $u \leq w.v$ . This context change is also reminiscent of the *context resurrection* operation of the literature on proof irrelevance (in the style of [Pfe01] for example).

We chose an equivalent but more conventional style where the context of sub-derivation does not change: instead of a judgment  $\Gamma \vdash \tau$  that becomes  $(\Gamma/u) \vdash \sigma$  when moving to a subterm of variance  $u$ , we define a judgment of the form  $\Gamma \vdash \tau : v$ , that evolves into  $\Gamma \vdash \sigma : (v.u)$ . The two styles are equally expressive: our judgment  $\Gamma \vdash \tau : v$  holds if and only if  $(\Gamma/v) \vdash \tau$  holds in Abel’s system—but we found that this one extends more naturally to checking decomposability, as will later be necessary. The inference rules for the judgment  $\Gamma \vdash \tau : v$  are defined on Figure 2.

**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  $\Gamma$ , the expression  $\tau$  varies along  $v$ ”. We also give a few formal results about this judgment.

<sup>8</sup>It turns out that this condition is not *necessary* and can be slightly weakened: we will discuss that later (3).

**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.

Proof:

*Soundness:*  $\Gamma \vdash \tau : v$  implies  $\llbracket \Gamma \vdash \tau : v \rrbracket$  By induction on the derivation. In the variable case this is direct. In the  $\bar{\sigma} \mathbf{t}$  case, for  $\bar{\rho}, \bar{\rho}'$  such that  $\bar{\rho} \prec_{\Gamma} \bar{\rho}'$ , we get  $\forall i, \sigma_i[\bar{\rho}] \prec_{v.w_i} \sigma_i[\bar{\rho}']$  by inductive hypothesis, which allows to conclude, by definition of variance composition, that  $(\bar{\sigma} \mathbf{t})[\bar{\rho}] \prec_v (\bar{\sigma} \mathbf{t})[\bar{\rho}']$ .

*Completeness:*  $\llbracket \Gamma \vdash \tau : v \rrbracket$  implies  $\Gamma \vdash \tau : v$  By induction on  $\tau$ ; in the variable case this is again direct. In the  $\bar{\sigma} \mathbf{t}$  case, given  $\bar{\rho} \prec_{\Gamma} \bar{\rho}'$  such that  $(\bar{\sigma} \mathbf{t})[\bar{\rho}] \prec_v (\bar{\sigma} \mathbf{t})[\bar{\rho}']$  we can deduce by inversion that for each variable  $\alpha_i$  of variance  $w_i$  in  $\tau[\bar{\alpha}]$  we have  $\sigma_i[\bar{\rho}] \prec_{v.w_i} \sigma_i[\bar{\rho}']$ , which allows us to inductively build the subderivations  $\Gamma \vdash \sigma_i : v.w_i$ . ■

**Lemma 2 (Monotonicity)** If  $\Gamma \vdash \tau : v$  is provable and  $\Gamma \leq \Gamma'$  then  $\Gamma' \vdash \tau : v$  is provable.

**Lemma 3** If  $\Gamma \vdash \tau : v$  and  $\Gamma' \vdash \tau : v$  both hold, then  $(\Gamma \vee \Gamma') \vdash \tau : v$  also holds.

**Corollary 1 (Principality)** For any type  $\tau$  and any variance  $v$ , there exists a minimal context  $\Delta$  such that  $\Delta \vdash \tau : v$  holds. That is, for any other context  $\Gamma$  such that  $\Gamma \vdash \tau : v$ , we have  $\Delta \leq \Gamma$ .

**Inversion of subtyping** We have mentioned in 2.1 the inversion properties of our subtyping relation. From  $\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}$  we can deduce subtyping relations on the type parameters  $\sigma_i, \sigma'_i$ . This can be generalized to any type expression  $\tau[\bar{\alpha}]$ :

**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  $\Gamma$  such that  $\bar{\sigma} \prec_{\Gamma} \bar{\sigma}'$ .

Proof: The reverse implication, is a direct application of the soundness of the variance judgment.

The direct implication is proved by induction on  $\tau[\bar{\alpha}]$ . The variable case is direct: if  $\alpha[\bar{\sigma}] \prec_v \alpha[\bar{\sigma}']$  holds then for  $\Gamma$  equal to  $(v\alpha)$  we indeed have  $v\alpha \vdash \alpha : v$  and  $\bar{\sigma} \prec_{\Gamma} \bar{\sigma}'$ .

In the  $\bar{\tau} \mathbf{t}$  case, we have that  $(\bar{\tau} \mathbf{t})[\bar{\sigma}] \prec_v (\bar{\tau} \mathbf{t})[\bar{\sigma}']$ . Suppose the variance of  $\bar{\alpha} \mathbf{t}$  is  $\bar{w}\bar{\alpha}$ : by inversion on the head type constructor  $\mathbf{t}$  we deduce that for each  $i$ ,  $\tau_i[\bar{\sigma}] \prec_{v.w_i} \tau_i[\bar{\sigma}']$ . Our induction hypothesis then gives us a family of contexts  $(\Gamma_i)_{i \in I}$  such that for each  $i$  we have  $\Gamma_i \vdash \tau_i : v.w_i$ . Furthermore,  $\bar{\sigma} \prec_{\Gamma_i} \bar{\sigma}'$  holds for all  $\Gamma_i$ , which means that  $\bar{\sigma} \prec_{\bigwedge_{i \in I} \Gamma_i} \bar{\sigma}'$ . Let's define  $\Gamma$  as  $\bigwedge_{i \in I} \Gamma_i$ . By construction we have  $\Gamma \geq \Gamma_i$ , so by monotonicity (Lemma 2) we have  $\Gamma \vdash \tau_i : v.w_i$  for each  $i$ . This allows us to conclude  $\Gamma \vdash \bar{\tau} \mathbf{t} : v$  as desired. ■

Note that this would work even for type constructors that are not  $v$ -closed: we are not comparing a  $\tau[\bar{\sigma}]$  to any type  $\tau'$ , but to a type  $\tau[\bar{\sigma}']$  sharing the same structure—the head constructors are always the same.

For any given pair  $\bar{\sigma}, \bar{\sigma}'$  such that  $\tau[\bar{\sigma}] \prec_v \tau[\bar{\sigma}']$  we can produce a context  $\Gamma$  such that  $\bar{\sigma} \prec_{\Gamma} \bar{\sigma}'$ . But is there a common context that would work for any pair? Indeed, that is the lowest possible context, the principal context  $\Gamma$  such that  $\Gamma \vdash \tau : v$ .



**Corollary 2 (Principal inversion)** *If  $\Delta$  is principal for  $\Delta \vdash \tau : v$ , then for any type sequences  $\bar{\sigma}$  and  $\bar{\sigma}'$ , the subtyping relation  $\tau[\bar{\sigma}] \prec_v \tau[\bar{\sigma}']$  implies  $\bar{\sigma} \prec_\Delta \bar{\sigma}'$ .*

**Proof:** Let  $\Delta$  be the principal context such that  $\Delta \vdash \tau : v$  holds. For any  $\bar{\sigma}, \bar{\sigma}'$  such that  $\tau[\bar{\sigma}] \prec_v \tau[\bar{\sigma}']$ , by inversion (Theorem 1) we have some  $\Gamma$  such that  $\Gamma \vdash \tau : v$  and  $\bar{\sigma} \prec_\Gamma \bar{\sigma}'$ . By definition of  $\Delta$ ,  $\Delta \leq \Gamma$  so  $\bar{\sigma} \prec_\Delta \bar{\sigma}'$  also holds. That is,  $\bar{\sigma} \prec_\Delta \bar{\sigma}'$  holds for any  $\bar{\sigma}, \bar{\sigma}'$  such that  $\tau[\bar{\sigma}] \prec_v \tau[\bar{\sigma}']$ . ■

**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 **type**  $\bar{v}\bar{\alpha} \mathbf{t} = (\mathbf{K}_c \text{ of } (\tau^c[\bar{\alpha}])_{c \in C})$  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  $\mathbf{K}$  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:

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

Instead of being just  $\alpha \text{ expr}$ , the codomain is now  $(\beta * \gamma) \text{ expr}$ . 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. The idea is that we will force again the codomain to be exactly  $\alpha \text{ expr}$ , but allow additional type equations such as  $\alpha = \beta * \gamma$  in this example:

$$\mid \text{Prod} : \forall \alpha. \forall \beta \gamma [\alpha = \beta * \gamma]. \beta \text{ expr} * \gamma \text{ expr} \rightarrow \alpha \text{ expr}$$

This restricted form justifies the name of *guarded* algebraic datatype. The  $\forall \beta \gamma [D]. \tau$  notation, coming from Simonet and Pottier, is a *constrained type scheme*:  $\beta, \gamma$  may only be instantiated with type parameters respecting the constraint  $D$ . Note that, as  $\beta$  and  $\gamma$  do not appear in the codomain anymore, we may equivalently change the outer universal into an existential on the left-hand side of the arrow:

$$| \text{Prod} : \forall \alpha. (\exists \beta \gamma [\alpha = \beta * \gamma]. \beta \text{ expr} * \gamma \text{ expr}) \rightarrow \alpha \text{ expr}$$

In the general case, a GADT definition for  $\bar{\alpha} \mathbf{t}$  is composed of a set of constructor declarations, each of the form:

$$| \text{K} : \forall \bar{\alpha}. (\exists \bar{\beta} [\bar{\alpha} = \bar{\sigma}[\bar{\beta}]]. \tau[\bar{\alpha}, \bar{\beta}]) \rightarrow \bar{\alpha} \mathbf{t}$$

or, reusing the classic notation,

$$| \text{K of } \exists \bar{\beta} [\bar{\alpha} = \bar{\sigma}[\bar{\beta}]]. \tau[\bar{\alpha}, \bar{\beta}]$$

Without loss of generality, we can conveniently assume that the variables  $\bar{\alpha}$  do not appear in the parameter type  $\tau$  anymore: if some  $\alpha_i$  appears in  $\tau$ , one may always pick a fresh existential variable  $\beta$ , add the constraint  $\alpha = \beta$  to  $D$ , and consider  $\tau[\beta/\alpha]$ . Let us re-express the introductory example in this form, that is,  $\text{K of } \exists \bar{\beta} [\bar{\alpha} = \bar{\sigma}[\bar{\beta}]]. \tau[\bar{\beta}]$ :

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

If all constraints between brackets are of the simple form  $\alpha_i = \beta_i$  (for distinct variables  $\alpha_i$  and  $\beta_i$ ), as for the constructor **Thunk**, then we have a constructor with existential types as described by Läufer and Odersky [OL92]. If furthermore there are no other existential variables than those equated with a type parameter, as in the **Val** case, we have an usual algebraic type constructor; of course the whole type is “simply algebraic” only if each of its constructors is algebraic.

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  $\text{type } \bar{v}\bar{\alpha} \mathbf{t}$  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  $\text{HMG}(X)$  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  $\bar{\alpha} \mathbf{t}$ , 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  $\tau_1, \tau_2, \bar{\sigma}, \bar{\sigma}'$  and distinct datatypes  $\bar{\alpha} \mathbf{t}, \bar{\alpha}' \mathbf{t}'$ , the types  $(\tau_1 \rightarrow \tau_2)$ ,  $\tau_1 * \tau_2$ ,  $\bar{\sigma} \mathbf{t}$  and  $\bar{\sigma}' \mathbf{t}'$  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 datatypes<sup>9</sup>: for each datatype  $\bar{\alpha} \mathbf{t}$  and all type vectors  $\bar{\sigma}$  and  $\bar{\sigma}'$  such that  $\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}$ , we must have  $(\exists \bar{\beta}[D[\bar{\sigma}]]\tau) \leq (\exists \bar{\beta}[D[\bar{\sigma}']]\tau)$  for each constructor  $K$  of  $\exists \bar{\beta}[D[\bar{\beta}, \bar{\alpha}]].\tau[\bar{\beta}]$ .

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  $\bar{\nu}\bar{\alpha} \mathbf{t}$  we are defining indeed has variance  $\bar{\nu}\bar{\alpha}$ , 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.

### 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.

**A remark on the non-completeness** Note that while the criterion REQ-SP is *sound*, it is not *complete*—even in the simple ADT case.

For a constructor  $K$  of  $\tau$  of  $\bar{\sigma} \mathbf{t}$ , the justification for the fact that, under the hypothesis  $\bar{\sigma} \mathbf{t} \leq \bar{\sigma}' \mathbf{t}$ , we should have  $\tau[\bar{\sigma}] \leq \tau[\bar{\sigma}']$  is the following: given a value  $v$  of type  $\tau[\bar{\sigma}']$ , we can build the value  $K v$  at type  $\bar{\sigma}' \mathbf{t}$ , and coerce it to  $\bar{\sigma} \mathbf{t}$ . We can then deconstruct this value by matching the constructor  $K$ , whose argument is of type  $\tau[\bar{\sigma}']$ . But this whole computation, (`match (K v :>  $\bar{\sigma}' \mathbf{t}$ ) with K x  $\rightarrow$  x`), reduces to  $v$ , so for value reduction to hold we need to also have  $v : \tau[\bar{\sigma}']$ .

For this whole argument to work we need a value at type  $\tau[\bar{\sigma}]$ . In fact, if the type  $\tau[\bar{\sigma}]$  is not inhabited, it can fail to satisfy REQ-SP and still be sound: this criterion is not complete. See the following example in a consistent system with an uninhabited type  $\perp$ :

```
type + $\alpha$  t =
  | T of int
  | Empty of  $\perp * (\alpha \rightarrow \text{bool})$ 
```

Despite  $\alpha$  occurring in a contravariant position in the dead `Empty` branch (which violates the soundness criterion of Simonet and Pottier), under the assumption that the  $\perp$  type really is uninhabited we know that this `Empty` constructor will never be used in closed code, and the contravariant occurrence will therefore never make a program “go wrong”. The definition is correct, yet rejected by REQ-SP, which is therefore incomplete.

<sup>9</sup>This is an extended version of the soundness requirement for algebraic datatypes: it is now formulated in terms guarded existential types  $\exists \bar{\beta}[D]\tau$  rather than simple argument types  $\tau$ .

Deciding type inhabitation in the general case is a very complex question, which is mostly orthogonal to the presence and design of GADT in the type system. There is, however, one clear interaction between the type inhabitation question and GADT. If a GADT  $\bar{\alpha} \mathbf{t}$  is instantiated with type variables  $\bar{\sigma}$  that satisfy none of the constraints  $D[\alpha]$  of its constructors  $K$  of  $\exists \bar{\beta}[\bar{D}].\tau$ , then we know that  $\bar{\sigma} \mathbf{t}$  is not inhabited. This is related to the idea of “domain information” that we discuss in the Future Work section (§5).

### 3.1 Expressing decomposability

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

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

We can substitute equalities and use the (assumed) covariance to simplify the subtyping constraint  $\sigma \text{ expr} \leq \sigma' \text{ expr}$  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  $\sigma'$  of a product  $\rho_1 * \rho_2$  is itself a product, *i.e.* is of the form  $\rho'_1 * \rho'_2$  for some  $\rho'_1$  and  $\rho'_2$ .

More generally, for a type  $\Gamma \vdash \sigma$  and a variance  $v$ , we are interested in a closure property of the form

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

Here, the context  $\Gamma$  represents the set of existential variables of the constructor ( $\beta$  and  $\gamma$  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  $\Gamma$  annotated with variances  $(+\beta, +\gamma)$ , and using the context ordering  $(\prec_\Gamma)$ . 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-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  $\Gamma$ , a type expression  $\sigma[\bar{\beta}]$  and two variances  $v$  and  $v'$ , we say that  $\sigma$  is *decomposable* under  $\Gamma$  from variance  $v$  to variance  $v'$ , which we write  $\Gamma \Vdash \sigma : v \rightsquigarrow v'$ , if the property

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

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  $\prec_\Gamma$  in the proposition  $\Gamma \Vdash \tau : v \rightsquigarrow v'$  and the anti-monotonicity of  $\prec_\Gamma$ , 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.*

Our final decomposability criterion, given below in Figure 3, requires both correct variances and a decomposability property, so it will be neither monotonous nor anti-monotonous with respect to the context argument.

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  $\beta, \gamma$ ). 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  $\gamma'$ . 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,  $\sigma'$  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  $\sigma'$  is equal to  $\sigma \text{ ref} * \sigma \text{ ref}$ , which is an instance<sup>10</sup> of  $\beta \text{ ref} * \beta \text{ ref}$ .

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 `type  $\alpha \text{ irr} = \text{int}$` ) is upward closed: if  $\sigma * (\sigma \text{ irr}) \leq \sigma'$ , then  $\sigma'$  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  $\sigma'$ , 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  $\sigma'$ . 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  $\beta_i$  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  $\sigma_i$ , 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  $\sigma'$  such that the two witnesses are incompatible—and similarly for one covariant and

<sup>10</sup>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.

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”<sup>11</sup> operation  $v_1 \lambda v_2$ , 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:

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

The following lemma uses zipping to merge together the results of the decomposition of several subterms  $(\sigma_i)_i$  into a “simultaneous decomposition”.

**Definition 4 (Simultaneous decomposition)** Given a context  $\Gamma$ , and families of type expressions  $(\sigma_i)_{i \in I}$  and variances  $(v_i)_{i \in I}$  and  $(v'_i)_{i \in I}$ , we define the following “simultaneous closure property”  $\Gamma \Vdash (T_i : v_i \rightsquigarrow v'_i)_{i \in I}$  defined as :

$$\forall(\bar{\rho} : \Gamma), \bar{\sigma}', (\forall i \in I, \sigma_i[\bar{\rho}] \prec_{v_i} \sigma'_i) \implies \exists(\bar{\rho}' : \Gamma), \bar{\rho} \prec_{\Gamma} \bar{\rho}' \wedge (\forall i \in I, \sigma_i[\bar{\rho}'] \prec_{v'_i} \sigma'_i)$$

**Lemma 5 (Soundness of zipping)** Suppose we have families of type expressions  $(T_i[\bar{\beta}])_{i \in I}$ , contexts  $(\Gamma_i)_{i \in I}$  and variances  $(v_i)_{i \in I}$  and  $(v'_i)_{i \in I}$  such that  $\lambda_{i \in I} \Gamma_i$  exists and for all  $i$  we have both  $\Gamma_i \vdash T_i : v_i$  and  $\Gamma_i \Vdash T_i : v_i \rightsquigarrow v'_i$ . Then, we have  $(\lambda_{i \in I} \Gamma_i) \Vdash (\sigma_i : v_i \rightsquigarrow v'_i)_{i \in I}$ .

**Proof:** Without loss of generality, we can consider that there are only two type expressions  $T_1[\bar{\beta}]$  and  $T_2[\bar{\beta}]$ , and that the free variables  $\bar{\beta}$  is reduced to a single variable  $\beta$ . Let  $w_1, w_2$  be the respective variances of  $\beta$  in  $\Gamma_1, \Gamma_2$ . We know that  $(w_1 \lambda w_2)$  exists and is equal to the variance  $w$  of the variable  $\beta$  in  $\Gamma_1 \lambda \Gamma_2$ .

Our further assumptions are  $\Gamma_i \vdash T_i : v_i$  (1) and  $\Gamma_i \Vdash T_i : v_i \rightsquigarrow v'_i$  (2) for  $i$  in  $\{1, 2\}$ . The expansion of (2) is:

$$\forall i \in I, \forall \rho, \sigma'_i, T_i[\rho] \prec_{v_i} \sigma'_i \implies \exists \rho', \rho \prec_{\Gamma_i} \rho' \wedge T_i[\rho'] \prec_{v'_i} \sigma'_i \quad (3)$$

Our goal is to prove  $\Gamma_1 \lambda \Gamma_2 \Vdash (\sigma_i : v_i \rightsquigarrow v'_i)_{i \in I}$ , which is equivalent to:

$$\forall \rho, \bar{\sigma}', (\forall i \in I, (T_i[\rho] \prec_{v_i} \sigma'_i)) \implies \exists \rho', \rho \prec_{\Gamma} \rho' \wedge \forall i \in I, T_i[\rho'] \prec_{v'_i} \sigma'_i \quad (4)$$

Assume given  $\rho$  and  $(\sigma'_1, \sigma'_2)$  such that  $T_1[\rho] \prec_{v_1} \sigma'_1$  (5) and  $T_2[\rho] \prec_{v_2} \sigma'_2$  (6). Applying (3) with  $i$  equal to 1 and (5) ensures the existence of a  $\rho'_1$  such that  $\rho \prec_{w_1} \rho'_1$  (7) and  $T_1[\rho'_1] \prec_{v_1} \sigma'_1$  (8). Similarly, there exists  $\rho'_2$  such that  $\rho \prec_{w_2} \rho'_2$  (9) and  $T_2[\rho'_2] \prec_{v_2} \sigma'_2$  (10). To establish (4), it remains to build a single  $\rho'$  that satisfies  $\rho \prec_w \rho'$  (11),  $T_1[\rho'] \prec_{v_1} \sigma'_1$  (12) and  $T_2[\rho'] \prec_{v_2} \sigma'_2$  (13), simultaneously. We reason by case analysis on  $w_1$  and  $w_2$  (restricted to the cases where the zip exists).

If both  $w_1$  and  $w_2$  are  $(\bowtie)$ , we take either  $\rho'_1$  or  $\rho'_2$  for: since  $w$  is  $\bowtie$ , we (11) trivially holds; Furthermore,  $T[\rho'] = T[\rho'_1] = T[\rho'_2]$  by irrelevance of  $v_i$  and (1); therefore (12) and (13) follow from (8) and (10).

If only one of the  $w_i$  is  $(\bowtie)$ , we'll suppose that it is  $w_1$ . We then take  $\rho_2$  for  $\rho'$ . Since  $\bowtie_1 \lambda w_2$  is  $w_2$ , (11) follows from (9); Furthermore,  $T[\rho'] = T[\rho'_1]$  by irrelevance of

<sup>11</sup>The idea of context merging and the term “zipping” are inspired by Montagu and Remy [MR09]

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

Figure 3: Syntactic decomposability

$v_1$  and (1) while  $T_2[\rho'] = T_2[\rho'_2]$  holds by construction; Hence, as in the previous case, (12) and (13) follow from (8) and (10).

Finally, if both  $w_1$  and  $w_2$  are  $(=)$ , then (7) and (9) implies  $\rho'_1 = \rho = \rho'_2$ . We take  $\rho$  for  $\rho'$  and all three conditions are obviously satisfied. ■

This lemma admits a kind of converse lemma stating completeness of zipping, that says that if  $\Gamma \Vdash (T_i : v_i \rightsquigarrow v'_i)_i$  holds, then  $\Gamma$  is indeed related to a zip of contexts  $(\Gamma_i)_i$  that pairwise decompose each of the  $(T_i)_i$ . However, the proof of completeness is more delicate and we prove it separately in §3.5.

### 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 rely on the zip soundness (Lemma 5) to merge sub-derivations into larger ones, so in addition to decomposability, the judgments simultaneously ensures that  $v$  is a correct variance for  $\tau$  under  $\Gamma$ . 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.

Rule SC-VAR is very similar to VC-VAR, except that the condition  $w \geq v$  is replaced by a stronger equality  $w = v$ . The reason why the variance-checking judgment has an inequality  $w \geq v$  is to make it monotonous in the environment—as requested by its corresponding semantics criterion (Definition 2). Therefore, the condition  $w \geq v$  is necessary for completeness—and admissible. On the contrary, the present judgment ensures, according the semantic criterion (Definition 5), that both the variance is correct (monotonous in the environment) and the type is decomposable, a property which is *anti-monotonic* in the environment (Lemma 4). Therefore, the semantics criterion  $\llbracket \Gamma \vdash \tau : v \Rightarrow v' \rrbracket$  is invariant in  $\Gamma$  and, correspondingly, the variable rule must use a strict equality.

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  $\Gamma_i$  that can be zipped into a unique environment  $\Gamma$ .

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 5 (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 6 (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.

**Proof:** The proof is by induction on the derivation  $\Gamma \vdash \tau : v \Rightarrow v'$  (1). Expanding  $\llbracket \Gamma \vdash \tau : v \Rightarrow v' \rrbracket$ , we must show both  $\llbracket \Gamma \vdash \tau : v \rrbracket$ , or equivalently  $\Gamma \vdash \tau : v$  (2) and  $\Gamma \Vdash \tau : v \rightsquigarrow v'$  (3), which itself expands to:

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

Let  $\bar{\rho}, \tau'$  be such that  $\tau[\bar{\rho}] \prec_v \tau'$ . We must exhibit a sequence  $\bar{\rho}'$  such that  $\bar{\rho} \prec_{\Gamma} \bar{\rho}'$  (4) and  $\tau[\bar{\rho}'] \prec_{v'} \tau'$  (5). Cases where the derivation of (1) ends with SC-TRIV and SC-VAR cases are direct: take  $\bar{\rho}$  and  $(\dots, \sigma', \dots)$  for  $\bar{\rho}'$ , respectively.

In the remaining cases, the derivation ends with Rule SC-CONSTR and  $\tau$  is of the form  $\bar{\sigma} \mathbf{t}$ .

- The  $v$ -closure assumption of the left premise ensures that  $\tau'$  is itself of the form  $\bar{\sigma}' \mathbf{t}$  for some sequence of closed types  $\bar{\sigma}'$ . By inversion on the variance  $\bar{w}\bar{\alpha}$  of the head constructor  $\mathbf{t}$ , we deduce  $\sigma_i[\bar{\rho}] \prec_{v.w_i} \sigma'_i$  for all  $i$  (6).
- The middle premise is the zipping assumption on the contexts  $\Gamma = \bigwedge_{i \in I} \Gamma_i$  (7).
- The right premises gives us subderivations  $\Gamma_i \vdash \sigma_i : v.w_i \Rightarrow v'.w_i$ . This implies  $\Gamma_i \vdash \sigma_i : v.w_i$ , for all  $i$ , which implies  $\Gamma \vdash \bar{\sigma} \mathbf{t} : v$ , *i.e.* (2). By induction hypothesis, this also implies  $\Gamma_i \Vdash \sigma_i : v.w_i$  (8) and  $\Gamma_i \Vdash \sigma_i : v.w_i \Rightarrow v'.w_i$  for all  $i$  (9).

We may now apply zip soundness (Lemma 5) with hypotheses (7), (8) and (9), which gives us the simultaneous decomposition  $\Gamma \Vdash (\sigma'_i : v.w_i \rightsquigarrow v'.w_i)_{i \in I}$ . Expanding this property (Definition 4), we may apply to (6) to get to get a witness  $\bar{\rho}'$  such that both  $\bar{\rho}' \prec_{\Gamma} \bar{\rho}$ , *i.e.* our first goal (4), and  $(\forall i \in I, \sigma_i[\bar{\rho}'] \prec_{v'.w_i} \sigma'_i)$ , which implies  $(\bar{\sigma} \mathbf{t})[\bar{\rho}'] \prec_{v'} \bar{\sigma}' \mathbf{t}$ , *i.e.* our second goal (5). ■

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. The proof requires several auxiliary lemmas; it is the subject of the next subsection.

### 3.5 Completeness of syntactic decomposability

We first show a few auxiliary results that will serve in the proof of zip completeness, and later, to reconnect our closure-checking criterion (Definition 5) with the full criterion of Simonet and Pottier (REQ-SP).

**Lemma 7 (Intermediate value)** Let  $\tau[\bar{\alpha}]$  be a type expression and  $\bar{\rho}_1, \bar{\rho}_2, \bar{\rho}_3$  three type families such that  $\tau[\bar{\rho}_1] \leq \tau[\bar{\rho}_2] \leq \tau[\bar{\rho}_3]$  and  $\bar{\rho}_1 \prec_{\Gamma} \bar{\rho}_3$  holds for some  $\Gamma$ . Then, there exists a type family  $\bar{\rho}'_2$  such that both  $\bar{\rho}_1 \prec_{\Gamma} \bar{\rho}'_2 \prec_{\Gamma} \bar{\rho}_3$  and  $\tau[\bar{\rho}'_2] = \tau[\bar{\rho}_2]$  hold.

**Proof:** We reuse the notations of the definition and assume  $\tau[\bar{\rho}_1] \leq \tau[\bar{\rho}_2] \leq \tau[\bar{\rho}_3]$  (1) and  $\bar{\rho}_1 \prec_{\Gamma} \bar{\rho}_3$  (2). We just have to exhibit  $\bar{\rho}'_2$  such that both  $\bar{\rho}_1 \prec_{\Gamma} \bar{\rho}'_2 \prec_{\Gamma} \bar{\rho}_3$  (3) and  $\tau[\bar{\rho}'_2] = \tau[\bar{\rho}_2]$  (4) hold. Let  $\Delta$  be the most general variance of  $\tau$ , *i.e.* the lowest context such that  $\Delta \vdash \tau : +$  (5) holds. By principal inversion (Corollary 2) applied to (1) twice thanks to (5), we have  $\bar{\rho}_1 \prec_{\Delta} \bar{\rho}_2 \prec_{\Delta} \bar{\rho}_3$  (6).



If  $\Delta \geq \Gamma$ , the result is immediate, as  $\bar{\rho}_1 \prec_{\Gamma} \bar{\rho}_2 \prec_{\Gamma} \bar{\rho}_3$  follows from (6) by anti-monotonicity and both (3) and (4) hold when we take  $\bar{\rho}_2$  for  $\bar{\rho}'_2$ . Otherwise, we reason on each variable of the context  $\Gamma$  independently. We may assume, *w.l.o.g.*, that  $\tau$  is defined over a single free variable  $\alpha$ , and  $\Gamma$  and  $\Delta$  are single variances  $v_{\Delta}, v_{\Gamma}$  with  $v_{\Delta} \not\prec v_{\Gamma}$ . We reason by case analysis on the possible variances for  $(v_{\Delta}, v_{\Gamma})$ , which are  $\{(-, =), (+, -), (-, +), (\bowtie, -)\}$ .

If  $v_{\Gamma}$  is  $(=)$ , the hypotheses (2) and (1) become  $\rho_1 = \rho_3$  and  $\tau[\rho_1] \leq \tau[\rho_2] \leq \tau[\rho_1]$ , which implies  $\tau[\rho_1] = \tau[\rho_2]$ . Thus, taking  $\rho_1$  for  $\rho'_2$  satisfies (3) and (4).

If  $v_{\Delta}$  is  $(\bowtie)$ , then by irrelevant of  $\bowtie$  and (5), we have  $\tau[\rho_1] = \tau[\rho_2]$ . Thus, taking  $\rho_1$  for  $\rho'_2$  satisfies (3) and (4), as above.

Finally, the cases  $(+, -)$  and  $(-, +)$  are symmetric and we will only work out the first one, *i.e.*  $v_{\Delta}$  is  $(+)$  and  $v_{\Gamma}$  is  $(-)$ . From  $\tau[\rho_1] \leq \tau[\rho_3]$  (which follows from (1) by transitivity) and (5), we have  $\rho_1 \prec_{v_{\Delta}} \rho_3$ , *i.e.*  $\rho_1 \leq \rho_3$ . Since the hypothesis (2) becomes  $\rho_1 \geq \rho_3$ , we have  $\rho_1 = \rho_3$ . Then, taking  $\rho_1$  for  $\rho'_2$ , (3) trivially holds while (4) follows from (1). ■

The next lemma connects the monotonicity of the variance-checking judgment (checking variance at a lower context provides more information, and is therefore harder) and the anti-monotonicity of the decomposability formula (decomposing to a higher context provides more information, and is therefore harder): for a fixed type expression, the contexts at which you can check variance are higher than the contexts at which you can decompose. This property, however, only holds for non-trivial decomposability results (otherwise any context can decompose): we must decompose from a  $v$  to a  $v'$  that do not verify  $v \geq v'$ , and no variable of the typing context must be irrelevant.

**Lemma 8** *Let  $\tau[\bar{\alpha}]$  be a type and  $v$  and  $v'$  be variances such that  $v \not\prec v'$ . If  $\Gamma \Vdash \tau : v \rightsquigarrow v'$  and  $\Delta$  is the most general context such that  $\Delta \vdash \tau : v$ , then, for each non-irrelevant variable  $\alpha$  of  $\Delta$ , we have  $\Gamma(\alpha) \leq \Delta(\alpha)$ .*

**Proof:** We show that the  $\Gamma$  is lower than the lowest possible  $\Delta$ , *i.e.* it is the most general context such that  $\Delta \vdash \tau : v$  holds. Without loss of generality, we consider the case where  $\tau$  has a single, non-irrelevant variable  $\alpha$ , and  $\Gamma$  and  $\Delta$  are singleton contexts over a single variance, respectively  $w_{\Gamma}$  and  $w_{\Delta}$ , with  $w_{\Delta} \neq (\bowtie)$ .

Therefore, we assume  $(w_{\Gamma}\alpha) \Vdash \tau : v \rightsquigarrow v'$  (1) and  $(w_{\Delta}\alpha) \vdash \tau : v$  (2). We prove that  $w_{\Gamma} \leq w_{\Delta}$ . We actually show that for any  $\rho_1, \rho_2$  such that  $\rho_1 \prec_{\Delta} \rho_2$  we also have  $\rho_1 \prec_{\Gamma} \rho_2$  (3).

From  $\rho_1 \prec_{\Gamma} \rho_2$ , we can deduce  $\tau[\rho_1] \prec_v \tau[\rho_2]$  (4). Applying (1), we get a  $\rho'$  such that  $\rho_1 \prec_{\Gamma} \rho'$  (5) and  $\tau[\rho'] \prec_{v'} \tau[\rho_2]$  (6). We then reason by case analysis on  $v \not\prec v'$ , considering the different cases  $\{(\bowtie, -), (+, -), (-, +), (-, =)\}$ .

If  $v$  is  $(\bowtie)$ , the most general  $w_{\Delta}$  is  $\bowtie$ , a case we explicitly ruled out: there is nothing to prove.

If  $v'$  is  $(=)$ , then (6) implies  $\tau[\rho'] = \tau[\rho_2]$ , and in particular  $\rho' = \rho_2$ ; our goal (3) follows from (5).

If  $(v, v')$  is  $(+, -)$  (we won't repeat the symmetric case  $(-, +)$ ), then (4) and (6) becomes  $\tau[\rho_1] \leq \tau[\rho_2]$  and  $\tau[\rho_2] \leq \tau[\rho']$ . Given (5), the intermediate value lemma (Lemma 7) ensures the existence of a  $\rho''$  such that  $\rho \prec_{\Gamma} \rho'' \prec_{\Gamma} \rho'$  and  $\tau[\rho''] = \tau[\rho_2]$ . From there, we deduce  $\rho'' = \rho_2$ , our goal (3) follows from (5), as in the previous case. ■

Finally, the following auxiliary lemmas will be useful in the proof of completeness.

**Lemma 9** *If the principal variance  $w$  such that  $(w\alpha) \vdash \tau[\alpha] : v$  holds is not the irrelevant variance  $\bowtie$ , then  $\tau[\rho_1] = \tau[\rho_2]$  implies  $\rho_1 = \rho_2$ .*

Proof: Whatever  $v$  is,  $\tau[\rho_1] = \tau[\rho_2]$  implies both  $\tau[\rho_1] \prec_v \tau[\rho_2]$  and its converse  $\tau[\rho_2] \prec_v \tau[\rho_1]$ . This holds in particular for the principal variance  $w$  such that  $(w\alpha) \vdash \tau[\alpha] : v$ . Moreover, by principal inversion (Corrolary 2) applied twice, we have both  $\rho_1 \prec_w \rho_2$  and  $\rho_2 \prec_w \rho_1$ . If  $w$  is distinct from  $\bowtie$  this implies  $\rho_1 = \rho_2$ . ■

**Lemma 10** *If  $\Gamma \Vdash \tau : v \rightsquigarrow (=)$  holds for some  $\Gamma$ , and  $\Delta$  is the most general context such that  $\Delta \vdash \tau : v$  holds, then  $\Delta \Vdash \tau : v \rightsquigarrow (=)$  also hold.*

Proof: Assume  $\Gamma \Vdash \tau : v \rightsquigarrow (=)$ , *i.e.*

$$\forall \bar{\rho}\bar{\sigma}', \tau[\bar{\rho}] \prec_v \bar{\sigma}' \implies \exists \bar{\rho}', \bar{\rho} \prec_{\Gamma} \bar{\rho}' \wedge \tau[\bar{\rho}'] = \bar{\sigma}' \quad (1)$$

Assume that  $\Delta$  is principal for  $\Delta \vdash \tau : v$  (2). We show  $\Delta \Vdash \tau : v \rightsquigarrow (=)$ , *i.e.*

$$\forall \bar{\rho}\bar{\sigma}', \tau[\bar{\rho}] \prec_v \bar{\sigma}' \implies \exists \bar{\rho}', \bar{\rho} \prec_{\Delta} \bar{\rho}' \wedge \tau[\bar{\rho}'] = \bar{\sigma}' \quad (3)$$

Let  $\bar{\rho}, \bar{\sigma}'$  be such that have  $\tau[\bar{\rho}] \prec_v \bar{\sigma}'$  (4). By (1), there exists  $\bar{\rho}'$  such that  $\tau[\bar{\rho}'] = \bar{\sigma}'$  (5). To prove (3), it only remains to prove that  $\bar{\rho} \prec_{\Delta} \bar{\rho}'$  (6). Given (5), the inequality (4) becomes  $\tau[\bar{\rho}] \prec_v \tau[\bar{\rho}']$  (7). Then (6) follows by principal inversion (Corrolary 2) applied to (7), given (2). ■

**Lemma 11** *Let  $\Delta$  be the most general context such that  $\Delta \vdash \tau : v$  holds. If  $v$  is  $(=)$ , then only variances  $(=)$  or  $(\bowtie)$  may appear in  $\Delta$ . If  $v$  is  $\bowtie$ , then only  $(\bowtie)$  may appear in  $\Delta$ .*

Proof: If  $v$  is  $(\bowtie)$ , the context  $\Gamma$  with all variances set to  $\bowtie$  satisfies  $\Gamma \vdash \tau : v$  (as  $\llbracket \Gamma \vdash \tau : v \rrbracket$  holds). By principality we have  $\Delta \leq \Gamma$ , so  $\Delta$  also has only irrelevant variances as  $(\bowtie)$  is the minimal variance.

If  $v$  is  $(=)$ , we handle each variable of the context independently, that is we can assume, *w.l.o.g.*, that  $\tau$  has only one variable  $\alpha$ . So  $\Delta$  is of the form  $(w\alpha)$  and we know that for any  $\rho, \rho'$  such that  $\rho \prec_w \rho'$  we have  $\tau[\rho] = \tau[\rho']$  (1). If  $w$  is not  $(\bowtie)$ , by lemma 9 applied to (1), we have  $\rho = \rho'$  for any  $\rho \prec_w \rho'$ , which means that  $w$  is  $(=)$ . Summing up, we have shown that  $w$  is either  $(\bowtie)$  or  $(=)$ . Reasoning similarly in the general case, any variance  $w$  of  $\Delta$  is either  $(\bowtie)$  or  $(=)$ . ■

We can now prove the converse of the zip soundness (Lemma 5) that is the core of the future proof of completeness of the decomposability judgment  $\Gamma \vdash \tau : v \Rightarrow (=)$ .

**Theorem 2 (Zip completeness)** *Given any context  $\Gamma$ , a family of type expressions  $(T_i[\bar{\sigma}])_{i \in I}$  and a family of variances  $(v_i)_{i \in I}$ , if the simultaneous decomposition  $\Gamma \Vdash (T_i : v_i \rightsquigarrow (=))_{i \in I}$  holds, then there exists a family of contexts  $(\Gamma_i)_{i \in I}$  such that  $\Gamma \leq \bigwedge_{i \in I} \Gamma_i$  and both  $\Gamma_i \vdash T_i : v_i$  and  $\Gamma_i \Vdash T_i : v_i \rightsquigarrow (=)$  hold for all  $i$ .*

*If furthermore  $\Gamma \vdash T_i : v_i$  holds for all  $i$ , then  $\bigwedge_{i \in I} \Gamma_i$  is precisely  $\Gamma$ .*

Proof: Let us assume the simultaneous decomposition  $\Gamma \Vdash (T_i : v_i \rightsquigarrow (=))_{i \in I}$ , which expands to:

$$\forall \bar{\sigma}', \bar{\rho}, (\forall i, T_i[\bar{\rho}] \prec_{v_i} \bar{\sigma}'_i) \implies \exists \bar{\rho}', \bar{\rho} \prec_{\Gamma} \bar{\rho}' \wedge (\forall i, T_i[\bar{\rho}'] = \bar{\sigma}'_i) \quad (1)$$

We construct a family of contexts  $(\Gamma_i)_{i \in I}$  such that the following holds:

$$\Gamma \leq \bigwedge_i \Gamma_i \text{ (2)} \quad \forall i, \Gamma_i \vdash T_i : v_i \text{ (3)} \quad \forall i, \Gamma_i \Vdash T_i : v_i \rightsquigarrow (=) \text{ (4)}$$

where (4) is equivalent to

$$\forall i, \forall \sigma'_i, \bar{\rho}, (T_i[\bar{\rho}] \prec_{v_i} \sigma'_i \implies \exists \bar{\rho}', \bar{\rho} \prec_{\Gamma_i} \bar{\rho}' \wedge T_i[\bar{\rho}'] = \sigma'_i) \text{ (5)}$$

The first step is to move from the entailment (1) of the form  $\forall \bar{\rho}, (\forall i \in I, \dots) \implies (\exists \bar{\rho}', (\forall i \in I, \dots))$  to the weaker form  $(\forall i \in I, \forall \bar{\rho}, \dots \implies \exists \bar{\rho}', \dots)$ , but closer to (5). More precisely, we show that

$$\forall i, \forall \sigma'_i, \bar{\rho}, (T_i[\bar{\rho}] \prec_{v_i} \sigma'_i \implies \exists \bar{\rho}', \bar{\rho} \prec_{\Gamma} \bar{\rho}' \wedge T_i[\bar{\rho}'] = \sigma'_i) \text{ (6)}$$

that is,  $\forall i, \Gamma \Vdash T_i : v_i \rightsquigarrow (=)$  (7). Let  $i, \sigma'_i$ , and  $\bar{\rho}$  be such that  $T_i[\bar{\rho}] \prec_{v_i} \sigma'_i$  (8). We show that there exists a  $\bar{\rho}'$  such that  $T_i[\bar{\rho}] = \sigma'_i$  (9) and  $\bar{\rho} \prec_{\Gamma} \bar{\rho}'$  (10). Let us extend our type  $\sigma'_i$  to a family  $\bar{\sigma}'$ , defined by taking  $\sigma'_j$  equal to  $T_j[\bar{\rho}]$  for  $j$  in  $I \setminus \{i\}$ . By construction, we have  $(\forall j \in I, T_j[\bar{\rho}] \prec_{v_j} \sigma'_j)$ . Therefore, we may apply (1) to get a  $\bar{\rho}'$  such that  $\bar{\rho} \prec_{\Gamma} \bar{\rho}'$ , *i.e.* our first goal (10), and  $(\forall j, T_j[\bar{\rho}'] = \sigma'_j)$ , which implies our second goal (9) when  $j$  is  $i$ . This proves (6).

We now prove that we can refine this to have  $\bar{\rho} \prec_{\Gamma_i} \bar{\rho}'$  for  $(\Gamma_i)_{i \in I}$  such that  $\Gamma \leq \bigwedge_i \Gamma_i$ .

Let  $\Delta_1$  and  $\Delta_2$  be the most general contexts such both that  $\Delta_1 \vdash T_1 : v_1$  and  $\Delta_2 \vdash T_2 : v_2$  hold (11). Let (2)', (3)', and (4)' be obtained by replacing  $\Gamma_i$ 's by  $\Delta_i$ 's in our three goals (2), (3), and (4). In fact (3)' is just (11). By Lemma 10 applied to (7) twice, given (11), we have both  $\Delta_1 \Vdash T_1 : v_1 \rightsquigarrow (=)$  and  $\Delta_2 \Vdash T_2 : v_2 \rightsquigarrow (=)$ , that is, (4).

Hence  $\Delta_1$  and  $\Delta_2$  are correct choices for  $\Gamma_1$  and  $\Gamma_2$  if they also satisfy the goal (2)', *i.e.*  $\Delta_1 \bigwedge \Delta_2 \geq \Gamma$ . We now study when remaining goal (2)' holds and, when it does not, propose a different choice for  $\Gamma_1$  and  $\Gamma_2$  that respect all three goals.

*W.l.o.g.*, we assume that  $I$  is reduced to  $\{1, 2\}$  and that there is only one free variable  $\beta$  in  $T_1, T_2$ . Since we focus on a single variable of the context we name  $w_1$  and  $w_2$  the variances of  $\beta$  in  $\Delta_1$  and  $\Delta_2$ , respectively. We now reason by case analysis on the variances  $w_1$  and  $w_2$ .

If both of them are  $\bowtie$ , we have  $\Delta_1 \bigwedge \Delta_2 = (\bowtie\beta)$ , so we do not necessarily have  $\Gamma \leq \Delta_1 \bigwedge \Delta_2$ . Instead, we make a different choice for  $\Gamma_2$ . Namely, we pick  $\Gamma$  for  $\Gamma_2$  and keep  $\Delta_1$  for  $\Gamma_1$ . As  $\Delta_2$  is  $\bowtie$  we have  $\Delta_2 \leq \Gamma_2$ , so by monotonicity of the variance checking judgment we have  $\Gamma_2 \vdash T_2 : v_2$  from (11), and we still have  $\Gamma_2 \Vdash T_2 : v_2 \rightsquigarrow (=)$  from (7). Hence (3) and (4) are reestablished. Finally, we have  $\Gamma_1 \bigwedge \Gamma_2 = \Gamma$ , so in particular  $\Gamma \leq \Gamma_1 \bigwedge \Gamma_2$ , *i.e.* (2).

If only one of the  $w_i$  is  $\bowtie$ , we may assume, *w.l.o.g.*, that it is  $w_1$ . Then  $\Delta_1 \bigwedge \Delta_2$  is  $\Delta_2$  and we only need to show that  $\Gamma \leq \Delta_2$  (12). From (7), we have  $\Gamma \Vdash T_2 : v_2 \rightsquigarrow (=)$ . We then make a case analysis on  $v_2$ : if  $v_2$  is not  $(=)$ , then by since  $\Delta_2$  is most general and  $w_2 \neq \bowtie$ , we may apply Lemma 8 to get (12); Otherwise,  $v_2$  is  $(=)$ ; Lemma 11 applied to (11) implies that  $\Delta_2$  is itself  $(=\beta)$ , and (12) trivially holds.

Finally, if none of the  $w_i$  is  $\bowtie$ , we first prove that they are both  $(=)$ . In fact, we only prove that  $w_1$  is  $(=)$  (13), as the other case follows by symmetry. To prove (13), we assume that  $\rho''$  be such that  $\rho \prec_{w_1} \rho''$  and we show that  $\rho = \rho''$  (14) holds. By (11), we have  $T_1[\rho] \prec_{v_1} T_1[\rho'']$ . By reflexivity, we have  $T_2[\rho] \prec_{v_2} T_2[\rho]$ . We can use those two inequalities to invoke our simultaneous decomposability hypothesis (1) with  $T_1[\rho'']$  for  $\sigma'_1$  and  $T_2[\rho]$  for  $\sigma'_2$  to get a  $\rho'$  such that both  $T_1[\rho'] = T_1[\rho'']$  and  $T_2[\rho'] = T_2[\rho]$  hold. By Lemma 9 applied with (11), this implies both  $\rho' = \rho''$  and  $\rho' = \rho$ , and therefore (14).

Therefore, the only remaining case is when  $w_1$  and  $w_2$  are both  $(=)$ . Then  $\Delta_1 \hat{\wedge} \Delta_2$  is  $(=\beta)$ , which is the highest single-variable context. So our goal (2)' trivially holds.

Of these several cases, one  $(\bowtie, \bowtie)$  has  $\Gamma_1 \hat{\wedge} \Gamma_2 = \Gamma$  directly, and in the others  $\Gamma_1$  and  $\Gamma_2$  were defined as the most general contexts such that  $\Gamma_1 \vdash T_1 : v_1$  and  $\Gamma_2 \vdash T_2 : v_2$ . If we add the further hypothesis that for each  $i$ ,  $\Gamma \vdash T_i : v_i$  holds, then by principality of the  $\Gamma_i$ , we have that  $\Gamma_i \leq \Gamma$  for each  $i$ . This implies that we have  $(\hat{\wedge}_{i \in I} \Gamma_i) \leq \Gamma$  (when it is defined,  $\hat{\wedge}$  coincides with the lowest upper bound  $\wedge$ ). By combination with (2), we get  $(\hat{\wedge}_{i \in I} \Gamma_i) = \Gamma$ . ■

**Lemma 12 (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.

**Proof:** Assume  $\llbracket \Gamma \vdash \tau : v \Rightarrow v' \rrbracket$  holds for  $v' \in \{=, \bowtie\}$ , i.e.  $\Gamma \vdash \tau : v$  (1) and  $\Gamma \Vdash \tau : v \rightsquigarrow v'$  (2), which expands to

$$\forall (\bar{\rho} : \Gamma), \tau', \tau[\bar{\rho}] \prec_v \tau' \implies \exists (\bar{\rho}' : \Gamma), \bar{\rho} \prec_\Gamma \bar{\rho}' \wedge \tau[\bar{\rho}'] \prec_{v'} \tau' \quad (3)$$

We show  $\Gamma \vdash \tau : v \Rightarrow v'$  (4) by structural induction on  $\tau$

If  $v \geq v'$  holds, then (4) directly follows from Rule SC-TRIV. This applies in particular when  $v' = \bowtie$ . Hence, we only need to consider the remaining cases where  $v'$  is  $(=)$  and  $v \not\geq v'$

We now reason by cases on  $\tau$ .

*Case  $\tau$  is a variable  $\alpha$ .* (3) becomes

$$\forall \rho, \tau', \rho \prec_v \tau' \implies \exists \rho', \rho \prec_\Gamma \rho' \wedge \rho' = \tau'$$

This means that if  $\rho \prec_v \tau'$  holds then  $\rho \prec_\Gamma \tau'$  also holds: the variance  $w\alpha \in \Gamma$  satisfies  $v \geq w$ . Since, the hypothesis (1) implies  $v \leq w$ , we have  $v = w$ . Therefore, (4) follows by Rule SC-VAR.

*Case  $\tau$  is of the form  $\bar{\sigma} \mathbf{t}$ .* By inversion, the derivation of (1) must end with rule VC-CONSTR, hence we have  $\Gamma \vdash \sigma_i : v.w_i$  (5) for each  $i \in I$  with  $\Gamma \vdash \mathbf{type} \bar{w}\alpha \mathbf{t}$  (6).

Let us show that  $\Gamma \Vdash (\sigma_i : v.w_i \rightsquigarrow =.w_i)_{i \in I}$  (7), i.e.

$$\forall \bar{\rho}, \bar{\sigma}', (\forall i \in I, \sigma_i[\bar{\rho}] \prec_{v_i} \sigma'_i) \implies \exists \bar{\rho}', \bar{\rho} \prec_\Gamma \bar{\rho}' \wedge (\forall i \in I, \sigma_i[\bar{\rho}'] \prec_{v'_i} \sigma'_i)$$

Let  $\bar{\rho}$  and  $\bar{\sigma}'$  be such that  $\sigma_i[\bar{\rho}] \prec_{v_i} \sigma'_i$  holds for all  $i$  in  $I$ . From this and (5), we have  $(\bar{\sigma} \mathbf{t})[\bar{\rho}] \prec_v \bar{\sigma}' \mathbf{t}$ . By application of (2), there exists  $\bar{\rho}'$  such that  $\bar{\rho} \prec_\Gamma \bar{\rho}'$  and  $(\bar{\sigma} \mathbf{t})[\bar{\rho}'] \prec_{=} \bar{\sigma}' \mathbf{t}$ . By inversion of subtyping, this implies  $\sigma_i[\bar{\rho}'] \prec_{=.w_i} \sigma'_i$ , for all  $i$  in  $I$ . This proves (7). We also note that the constructor  $\mathbf{t}$  is  $v$ -closed (8).

To prove our goal (4), we construct a family  $(\Gamma_i)_{i \in I}$  of contexts that satisfies  $\hat{\wedge}_{i \in I} \Gamma_i = \Gamma$  (9) and subderivations  $\Gamma_i \vdash \sigma_i : v.w_i \Rightarrow =.w_i$  (10), since then the conclusion (4) follows by an application of rule SC-CONSTR with (6), (9), and (10).

We will handle separately the arguments  $\sigma_j$  that are irrelevant, i.e. when  $w_j$  is  $\bowtie$ , from the rest. Let  $I_\bowtie$  be the set of indices with irrelevant variances and  $I_\bowtie$  the others.

For any  $i \in I_\bowtie$ ,  $v.w_i$  and  $=.w_i$  are both  $\bowtie$  so the condition (10), which becomes  $\Gamma_i \vdash \sigma_i : \bowtie \Rightarrow \bowtie$ , is void of content ( $\sigma_i[\bar{\rho}] \prec_\bowtie \sigma'_i \implies \exists \bar{\rho}', \sigma_i[\bar{\rho}'] \prec_\bowtie \sigma'_i$  is always true). More precisely, let  $\Gamma_i$  be the irrelevant context having only irrelevant variances. Then (10) follows by Rule SC-TRIV.

Since the decomposability constraints for  $i \in I_\bowtie$  such that  $w_i = \bowtie$  are trivial, (7) is equivalent to  $\Gamma \Vdash (\sigma_i : v.w_i \rightsquigarrow =.w_i)_{i \in I_\bowtie}$  (11).

For each  $i \in I_{\mathfrak{M}}$ ,  $(=.w_i)$  equals  $(=)$ , so (11) becomes  $\Gamma \Vdash (\sigma_i : v.w_i \rightsquigarrow (=))_{i \in I_{\mathfrak{M}}}$ . By zip completeness (Theorem 2), there is a family  $(\Gamma_i)_{i \in I_{\mathfrak{M}}}$  such that  $\Gamma \leq \lambda_{i \in I_{\mathfrak{M}}} \Gamma_i$  and both  $\Gamma_i \vdash \sigma_i : v.w_i$  and  $\Gamma_i \Vdash \sigma_i : v.w_i \rightsquigarrow (=)$ , i.e.  $\llbracket \Gamma_i \vdash \sigma_i : v.w_i \Rightarrow (=) \rrbracket$  (12) hold for any  $i \in I_{\mathfrak{M}}$ . Furthermore, since we also have (5), we can strengthen our result into  $\lambda_{i \in I_{\mathfrak{M}}} \Gamma_i = \Gamma$ . By induction hypothesis applied to (12), we have (10) for  $i \in I_{\mathfrak{M}}$ .

We have two families of contexts over domains  $I_{\mathfrak{M}}$  and  $I_{\mathfrak{M}}$  that partition  $I$ ; we can union them in a family  $(\Gamma_i)_{i \in I}$  that has subderivations  $\forall i \in I, \Gamma_i \vdash \sigma_i : v.w_i \Rightarrow v'.w_i$ . As the contexts in  $I_{\mathfrak{M}}$  are all irrelevant, they are neutral for the zipping operation:  $\lambda_{i \in I} \Gamma_i$  is equal to  $\lambda_{i \in I_{\mathfrak{M}}} \Gamma_i$ , that is  $\Gamma$ . This proves (9) while (10) has already been proved separately for  $i \in I_{\mathfrak{M}}$  and  $i \in I_{\mathfrak{M}}$ . ■

**Remark 1 (Note (8))** *The head constructor  $\mathfrak{t}$  is closed in our system with atomic subtyping, but the situation is in fact a bit stronger than that: the statement of  $v$ -closure of  $\bar{\alpha} \mathfrak{t}$  can be formulated in term of decomposability  $\Gamma \Vdash \bar{\alpha} \mathfrak{t} : v \rightsquigarrow (=)$ . It is very close from our decomposability hypothesis  $\Gamma \Vdash \bar{\sigma} \mathfrak{t} : v \rightsquigarrow (=)$ , but uses variables  $\bar{\alpha}$  instead of full type expressions  $\bar{\sigma}$ . We conjecture that the decomposability hypothesis (with  $=$  on the right) implies  $v$ -closure in a much larger set of subtyping systems that just atomic subtyping: it suffices that the subtyping relation is defined only in term of head constructors.*

### 3.6 Back to the correctness criterion

Remember the correctness criterion REQ-SP of Simonet and Pottier:

$$\forall \bar{\sigma}, \bar{\sigma}', \bar{\rho}, \quad (\bar{\sigma} \mathfrak{t} \leq \bar{\sigma}' \mathfrak{t} \wedge D[\bar{\sigma}, \bar{\rho}] \implies \exists \bar{\rho}', D[\bar{\sigma}', \bar{\rho}'] \wedge \tau[\bar{\rho}] \leq \tau[\bar{\rho}']) \quad (1)$$

We now show how the closure judgment  $\Gamma \vdash \tau : v \Rightarrow v'$  can be used to verify that this criterion holds: we will express this criterion in an equivalent form that uses the interpretation of our judgments.

The first step is to rewrite the property  $\bar{\sigma} \mathfrak{t} \leq \bar{\sigma}' \mathfrak{t}$  using the variance annotation  $\bar{v}\bar{\alpha}$  of  $\mathfrak{t}$ . Again, we are taking the variance annotation for the datatype  $\mathfrak{t}$  as granted (this is why we can use it in this reasoning step), and checking that the definitions of the constructors of  $\mathfrak{t}$  are sound with respect to this annotation.

$$\forall \bar{\sigma}, \bar{\sigma}', \bar{\rho}, \quad ((\forall i, \sigma_i \prec_{v_i} \sigma'_i) \wedge D[\bar{\sigma}, \bar{\rho}] \implies \exists \bar{\rho}', D[\bar{\sigma}', \bar{\rho}'] \wedge \tau[\bar{\rho}] \leq \tau[\bar{\rho}']) \quad (2)$$

Since, the constraint  $D[\bar{\alpha}, \bar{\beta}]$  is a set of equalities of the form  $\alpha_i = T_i[\bar{\beta}]$  (where  $T_i$  is a type), (2) is actually:

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

Substituting the equalities and, in particular, removing the quantification on the  $\bar{\sigma}$ , which are fully determined by the equality constraints  $\bar{\sigma} = \bar{T}[\bar{\rho}]$ , we get:

$$\forall \bar{\sigma}', \bar{\rho}, \quad (\forall i, T_i[\bar{\rho}] \prec_{v_i} \sigma'_i) \implies \exists \bar{\rho}', (\forall i, T_i[\bar{\rho}'] = \sigma'_i) \wedge \tau[\bar{\rho}] \leq \tau[\bar{\rho}'] \quad (3)$$

By inversion (Theorem 1), we may replace the goal  $\tau[\bar{\rho}] \leq \tau[\bar{\rho}']$  by the formula  $\exists \Gamma, (\Gamma \vdash \tau : +) \wedge (\bar{\rho} \prec_{\Gamma} \bar{\rho}')$ . Moreover, since  $\Gamma \vdash \tau : +$  always for for some  $\Gamma$ , we may move this quantification in front. Hence, (3) is equivalent to:

$$\exists \Gamma, \bigwedge \left\{ \begin{array}{l} \Gamma \vdash \tau : + \\ \forall \bar{\sigma}', \bar{\rho}, (\forall i, T_i[\bar{\rho}] \prec_{v_i} \sigma'_i) \implies \exists \bar{\rho}', (\forall i, T_i[\bar{\rho}'] = \sigma'_i) \wedge \bar{\rho} \prec_{\Gamma} \bar{\rho}' \end{array} \right. \quad (4)$$

We may recognize in second clause the simultaneous decomposability judgment (Definition 3)  $\Gamma \vdash (T_i : v_i \rightsquigarrow =)_{i \in I}$ . Hence, (4) is in fact:

$$\exists \Gamma, \Gamma \vdash \tau : + \wedge \Gamma \vdash (T_i : v_i \rightsquigarrow =)_{i \in I} \quad (5)$$

Then, comes the delicate step of this series of equivalent rewriting:

$$\exists \Gamma, (\Gamma_i)_{i \in I}, \bigwedge \left\{ \begin{array}{l} \Gamma \vdash \tau : + \\ \Gamma = \bigwedge_{i \in I} \Gamma_i \wedge \forall i, (\Gamma_i \vdash T_i : v_i \wedge \Gamma_i \vdash T_i : v_i \rightsquigarrow =) \end{array} \right. \quad (6)$$

The reverse implication from (6) to (5) is the zip soundness (Lemma 5).

The direct implication, from (5) to (6) is more involved: let  $\Gamma_0$  be such that  $\Gamma_0 \vdash \tau : +$ . By zip completeness (Theorem 2), with the hypotheses of (4), there exists a family  $(\Gamma_i)_{i \in I}$  satisfying the typing, zipping and decomposability of second line of (6) with  $\Gamma_0 \leq \bigwedge_{i \in I} \Gamma_i$ . We take  $\bigwedge_{i \in I} \Gamma_i$  for  $\Gamma$ . Then, from  $\Gamma_0 \leq \Gamma$  we get  $\Gamma \vdash \tau : +$  by monotonicity (Lemma 2).

As a last step, the last conjuncts of (6) are equivalent to  $\forall i, \Gamma_i \vdash T_i : v_i \Rightarrow (=)$  by interpretation of syntactic decomposability (Definition 5) and soundness and completeness of zipping (lemmas 6 and 12). Therefore, (6) is equivalent to:

$$\exists \Gamma, (\Gamma_i)_{i \in I}, \Gamma \vdash \tau : (+) \wedge \Gamma = \bigwedge_{i \in I} \Gamma_i \wedge \forall i \in I, \Gamma_i \vdash T_i : v_i \Rightarrow (=) \quad (7)$$

which is our final criterion.

**Pragmatic evaluation of this 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 are only concerned with the  $\Gamma \vdash \tau : +$  part.

To compute the contexts  $\Gamma$  and  $(\Gamma_i)_{i \in I}$  existentially quantified in this formula, one can use a variant of our syntactic judgments where the environment  $\Gamma$  is not an input, but an output of the judgment; in fact, one should return for each variable  $\alpha$  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 here. 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 possibles variances for  $\tau$ , 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 Upward and downward closure in a ML type system

In the type system we have used so far, *all types* are both upward and downward-closed. Indeed, thanks to the simplicity of our subtyping relation, we have a very strong inversion principle: two ground types in a subtyping relation necessarily have exactly the same structure. We have therefore completely determined a sound variance check for a simple type system with GADT.

This simple resolution, 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  $\sigma$  we have  $\perp \leq \sigma$ , downward-closure fails for most types. For example, products are no longer downward-closed:  $\Gamma \vdash \sigma * \tau \geq \perp$  does not imply 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  $\top$  type<sup>12</sup>. However, object types and polymorphic variant have subtyping, so they are, in general, neither upward nor downward-closed. Finally, subtyping is also used in private type definitions, that 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 reasonably large set of use cases.

## 4.2 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  $\tau$`  would be rejected by the type-checker if  $\tau$  was itself declared `downward-closed`.

Such “closure specifications” are part of the semantic properties of a type and would, as such, sometimes need to be exposed through module boundaries. It is important that the specification language for abstract types allow to say that a type is upward-closed (respectively downward-closed). These new ways to classify types raise some software engineering questions. When is it desirable to define types as `upward-closed`? The user must balance its ability to define semi-abstract version of the type against its use in a GADT—and potentially other type-system features that would make use of negative reasoning on the subtyping relation. We do not yet know how to answer this question and believe that more practice is necessary to get a clearer picture of the trade-off involved.

## 4.3 Subtyping constraints and variance assignment

We will now revisit our previous example, using the guarded existential notation:

<sup>12</sup>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.

```

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  $\exists\beta\gamma[\alpha = \beta * \gamma].\beta \text{ expr} * \gamma \text{ expr}$ 

```

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  $\exists\beta\gamma[\alpha \geq \gamma].\beta \text{ expr} * (\beta \rightarrow \gamma)$ 
  | Prod of  $\exists\beta\gamma[\alpha \geq \beta * \gamma].\beta \text{ expr} * \gamma \text{ expr}$ 

```

We have turned each equality constraint  $\alpha = T[\bar{\beta}]$  into a subtyping constraint  $\alpha \geq T[\bar{\beta}]$ . For a type  $\alpha'$  such that  $\alpha \leq \alpha'$ , we get by transitivity that  $\alpha' \geq T[\bar{\beta}]$ . This means that  $\alpha \text{ expr}$  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<sup>13</sup>: 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  $\text{eval} : \alpha \text{ expr} \rightarrow \alpha$ , because the constraints  $\alpha \geq \tau$  are in the right direction to get an  $\alpha$  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) :>  $\alpha$ )

```

We conjecture that moving from an equality constraint on the GADT type parameters to a subtyping constraint (bigger than, or smaller than, according to the desired variance for the parameter) is often unproblematic in practice. In the examples we have studied, such a change did not stop functions from type-checking—we only needed to add some explicit coercions.

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. This is related to practical implementation questions, as languages based on inference by unification tend to favor equality over subtyping, bidirectional coercions over unidirectional ones. 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

<sup>13</sup> Note that the formal proofs of the precedent section were, in some cases, specialized to the equality constraint. More precisely, our decomposability criterion is still sound when extended to arbitrary subtyping constraints, but its completeness is unknown and left to future work.



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.

**On the relaxed value restriction** Regarding the relaxed value restriction, which is our initial practical motivation to investigate variance in presence of GADT, there is also future work to be done to verify that it is indeed compatible with this refined notion of variance. While the syntactic proof of soundness of the relaxation doesn't involve subtyping directly, the "informal justification" for value restriction uses the admissibility of a global bottom type  $\perp$  to generalize a covariant unification variable; in presence of downward-closed type, there is no such general  $\perp$  type (only one for non-downward-closed types). We conjecture that the relaxed value restriction is still sound in this case, because the covariance criterion is really used to rule out mutable state rather than subtype from a  $\perp$  type; but it will be necessary to study the relaxation justification in more details to formally establish this result.

**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 §4.3, 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.

On the implementation side, we suspect that adding subtyping constraints to a type system that already supports GADT and private types should not require large engineering efforts (in particular, it does not imply supporting the most general forms of bounded polymorphism). Matching on a GADT  $\alpha$  `t` already introduces local type equalities of the form  $\alpha = T[\bar{\beta}]$  in pattern matching clauses. Jacques Garrigue suggested that adding an equality of the form  $\alpha = \text{private } T[\bar{\beta}]$  should correspond to GADT equations of the form  $\alpha \leq T[\bar{\beta}]$ , and lower bounds could be represented using the dual notion of `invisible` types. Regardless of implementation difficulties, in a system with only explicit subtyping coercion, such subtyping constraints would still require more user annotations.

**Mathematical structures for variance studies** There has been work on more structured presentation of GADT as part of a categorical framework ([GJ08] and [HF11]). This is orthogonal to the question of variance and subtyping, but it may be interesting to re-frame the current result in this framework.

Parametrized types with variance can also be seen as a sub-field of order theory with very partial orders and functions with strong monotonicity properties. Finally, we have been surprised to find that geometric intuitions were often useful to direct our formal developments. It is possible that existing work in these fields would allow us to streamline the proofs, which currently are rather low-level and tedious.

**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, the equality type  $(\alpha, \beta) \text{ eq}$  cannot be accurately described by considering variation of each of its parameters independently. We would like to say that  $(\alpha, \beta) \text{ eq} \leq (\alpha', \beta') \text{ eq}$  holds as soon as  $\alpha = \beta$  and  $\alpha' = \beta'$ . With the simple notion of variance we currently have, all we can soundly say about  $\text{eq}$  is that it must be invariant in both its parameters—which is considerably weaker. In particular, the well-known trick of “factoring out” GADT by using the  $\text{eq}$  type in place of equality constraint does not preserve variances: equality constraints allow fine-grained variance considerations based on upward or downward-closure, while the equality type instantly makes its parameters invariant.

We think it would be possible to regain some “completeness”, and in particular re-enable factoring by  $\text{eq}$ , 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  $\bar{\sigma} \text{ t} \leq \bar{\sigma}' \text{ t}$  to only cases where  $\bar{\sigma} \text{ t}$  and  $\bar{\sigma}' \text{ t}$  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  $(\alpha, \beta) \text{ eq}$  is restricted by the constraint  $\alpha = \beta$ , we could soundly declare the variance  $(\bowtie \alpha, \bowtie \beta) \text{ eq}$  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.
- [GJ08] Neil Ghani and Patricia Johann. Foundations for structured programming with gadts. In *Proceedings of Principles and Programming Languages (POPL), 2008*, pages 297–308, 2008.
- [HF11] Makoto Hamana and Marcelo Fiore. A foundation for gadts and inductive families: dependent polynomial functor approach. In *Proceedings of the seventh ACM SIGPLAN workshop on Generic programming*, 2011. URL: <http://www.cs.gunma-u.ac.jp/~hamana/Papers/dep.pdf>.
- [Kis] Oleg Kiselyov. Typed tagless interpretations and typed compilation. URL: <http://okmij.org/ftp/tagless-final/index.html>.
- [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>.
- [OL92] Martin Odersky and Konstantin Läufer. An extension of ML with first-class abstract types. In *ACM Workshop on ML and its Applications*, pages 78–91, June 1992. URL: <http://www.cs.luc.edu/lauffer/papers/ml92.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.

---

## Contents

<b>1</b>	<b>Motivation</b>	<b>3</b>
<b>2</b>	<b>A formal setting</b>	<b>9</b>
2.1	Atomic subtyping . . . . .	9
2.2	The algebra of variances . . . . .	10
2.3	Variance assignment in ADTs . . . . .	11
2.4	Variance annotations in GADT . . . . .	14
<b>3</b>	<b>Checking variances of GADT</b>	<b>16</b>
3.1	Expressing decomposability . . . . .	17
3.2	Variable occurrences . . . . .	18
3.3	Context zipping . . . . .	18
3.4	Syntactic decomposability . . . . .	20
3.5	Completeness of syntactic decomposability . . . . .	21
3.6	Back to the correctness criterion . . . . .	26
<b>4</b>	<b>Closed-world vs. open-world subtyping</b>	<b>27</b>
4.1	Upward and downward closure in a ML type system . . . . .	27
4.2	A better control on upward and downward-closure . . . . .	28
4.3	Subtyping constraints and variance assignment . . . . .	28
<b>5</b>	<b>Future Work</b>	<b>30</b>



**RESEARCH CENTRE  
PARIS – ROCQUENCOURT**

Domaine de Voluceau, - Rocquencourt  
B.P. 105 - 78153 Le Chesnay Cedex

Publisher  
Inria  
Domaine de Voluceau - Rocquencourt  
BP 105 - 78153 Le Chesnay Cedex  
[inria.fr](http://inria.fr)

ISSN 0249-0803