

# InvarGenT: Implementation

BY ŁUKASZ STAFINIĄK

Institute of Computer Science  
University of Wrocław

## Abstract

InvarGenT is a proof-of-concept system for invariant generation by full type inference with Guarded Algebraic Data Types and existential types encoded as automatically generated GADTs. This implementation documentation focuses on source code, refers to separate technical reports on theory and algorithms.

## 1 Data Structures and Concrete Syntax

Following [1], we have the following nodes in the abstract syntax of patterns and expressions:

- p-Empty.** 0: Pattern that never matches. Concrete syntax: `!`. Constructor: `Zero`.
- p-Wild.** 1: Pattern that always matches. Concrete syntax: `_`. Constructor: `One`.
- p-And.**  $p_1 \wedge p_2$ : Conjunctive pattern. Concrete syntax: e.g. `p1 as p2`. Constructor: `PAnd`.
- p-Var.**  $x$ : Pattern variable. Concrete syntax: lower-case identifier e.g. `x`. Constructor: `PVar`.
- p-Cstr.**  $Kp_1 \dots p_n$ : Constructor pattern. Concrete syntax: e.g. `K (p1, p2)`. Constructor: `PCons`.
- Var.**  $x$ : Variable. Concrete syntax: lower-case identifier e.g. `x`. Constructor: `Var`. External functions are represented as variables in global environment.
- Cstr.**  $Ke_1 \dots e_n$ : Constructor expression. Concrete syntax: e.g. `K (e1, e2)`. Constructor: `Cons`.
- App.**  $e_1 e_2$ : Application. Concrete syntax: e.g. `x y`. Constructor: `App`.
- LetRec.** `letrec  $x = e_1$  in  $e_2$` : Recursive definition. Concrete syntax: e.g. `let rec f = function ... in ...`. Constructor: `Letrec`.
- Abs.**  $\lambda(c_1 \dots c_n)$ : Function defined by cases. Concrete syntax: for single branching via `fun` keyword, e.g. `fun x y -> f x y` translates as  $\lambda(x.\lambda(y.(fx) y))$ ; for multiple branching via `match` keyword, e.g. `match e with ...` translates as  $\lambda(\dots)e$ . Constructor: `Lam`.
- Clause.**  $p.e$ : Branch of pattern matching. Concrete syntax: e.g. `p -> e`.
- CstrIntro.** Does not figure in neither concrete nor abstract syntax. Scope of existential types is thought to retroactively cover the whole program.
- ExCases.**  $\lambda[K](p_1.e_1 \dots p_n.e_n)$ : Function defined by cases and abstracting over the type of result. Concrete syntax: `function` and `ematch` keywords – e.g. `function Nil -> ... | Cons (x,xs) -> ...; ematch l with ...`. Parsing introduces a fresh identifier for  $K$ . Constructor: `ExLam`.
- ExLetIn.** `let  $p = e_1$  in  $e_2$` : Elimination of existentially quantified type. Concrete syntax: e.g. `let v = f e ... in ...`. Constructor: `Letin`.

We also have one sort-specific type of expression, numerals.

For type and formula connectives, we have ASCII and unicode syntactic variants (the difference is only in lexer). Quantified variables can be space or comma separated. The table below is analogous to information for expressions above. Existential type construct introduces a fresh identifier for  $K$ . The abstract syntax of types is not sort-safe, but type variables carry sorts which are inferred after parsing. Existential type occurrence in user code introduces a fresh identifier, a new type constructor in global environment `newtype_env`, and a new value constructor in global environment `newcons_env` – the value constructor purpose is to store the content of the existential type, it is not used in the program.

type variable	$x$	<code>x</code>		<code>TVar</code>
type constructor	<code>List</code>	<code>List</code>		<code>TCons(CNamed...)</code>
number (type)	<code>7</code>	<code>7</code>		<code>NCst</code>
numeral (expr.)	<code>7</code>	<code>7</code>		<code>Num</code>
numerical sum (type)	$a + b$	<code>a+b</code>		<code>Nadd</code>
existential type	$\exists \alpha \beta [a \leq \beta]. \tau$	<code>ex a b [a&lt;=b] . t</code>	$\exists a, b [a \leq b] . t$	<code>TCons(Extype...)</code>
type sort	$s_{ty}$	<code>type</code>		<code>Type_sort</code>
number sort	$s_R$	<code>num</code>		<code>Num_sort</code>
function type	$\tau_1 \rightarrow \tau_2$	<code>t1 -&gt; t2</code>	$t1 \rightarrow t2$	<code>Fun</code>
equation	$a \doteq b$	<code>a = b</code>		<code>Eqty</code>
inequation	$a \leq b$	<code>a &lt;= b</code>	$a \leq b$	<code>Leq</code>
conjunction	$\varphi_1 \wedge \varphi_2$	<code>a=b &amp;&amp; b=a</code>	$a=b \wedge b=a$	built-in lists

Toplevel expressions (corresponding to structure items in OCaml) introduce types, type and value constructors, global variables with given type (external names) or inferred type (definitions).

type constructor	<code>newtype List : type * num</code>	<code>TypConstr</code>
value constructor	<code>newcons Cons : all n a. a * List(a,n) --&gt; List(a,n+1)</code>	<code>ValConstr</code>
	<code>newcons Cons : <math>\forall n, a. a * List(a,n) \rightarrow List(a,n+1)</math></code>	
declaration	<code>external filter : <math>\forall n, a. List(a,n) \rightarrow \exists k [k \leq n]. List(a,k)</math></code>	<code>PrimVal</code>
rec. definition	<code>let rec f =...</code>	<code>LetRecVal</code>
non-rec. definition	<code>let v =...</code>	<code>LetVal</code>

For simplicity of theory and implementation, mutual non-nested recursion and or-patterns are not provided. For mutual recursion, nest one recursive definition inside another.

## 2 Generating and Normalizing Formulas

We inject the existential type and value constructors during parsing for user-provided existential types, and during constraint generation for inferred existential types, into the list of toplevel items, which allows to follow [1] despite removing `extype` construct from the language. It also facilitates exporting inference results as OCaml source code.

Functions `constr_gen_pat` and `envfrag_gen_pat` compute formulas according to table 2 in [1], and `constr_gen_expr` computes table 3. We preserve the FOL language presentation in the type `cnstrnt`, only limiting the expressivity in ways not requiring any preprocessing. The toplevel definitions (from type `struct_item`) `LetRecVal` and `LetVal` are processed by `constr_gen_letrec` and `constr_gen_let` respectively. They are analogous to `Letrec` and `Letin` or a `Lam` clause. We do not cover toplevel definitions in our formalism (without even a rudimentary module system, the toplevel is a matter of pragmatics rather than semantics).

Toplevel definitions (and in future, structure items) are intended as boundaries for constraint solving. This way the programmer can decompose functions that could be too complex for the solver. `LetRecVal` only binds a single identifier, while `LetVal` binds variables in a pattern. To preserve the flexibility of expression-level pattern matching, `LetVal` has to pack the constraints  $\llbracket \Sigma \vdash p \uparrow \alpha \rrbracket$  which the pattern makes available, into existential types. Each pattern variable is a separate entry to the global environment, therefore the connection between them is lost.

The formalism (in interests of parsimony) requires that only values of existential types be bound using `Letin` syntax. The implementation is enhanced in this regard: if the normalization step cannot determine which existential type is being eliminated, the constraint is replaced by one that would be generated for a pattern matching branch. This recovers the common use of the `let...in` syntax, with exception of polymorphic `let` cases, where `let rec` still needs to be used.

In the formalism, we use  $\mathcal{E} = \{\varepsilon_K, \chi_K | K :: \forall \alpha \gamma [\chi_K(\alpha, \gamma)]. \gamma \rightarrow \varepsilon_K(\alpha) \in \Sigma\}$  for brevity, as if all existential types  $\varepsilon_K(\alpha)$  were related with a predicate variable  $\chi_K(\alpha, \gamma)$ . In the implementation, we have user-defined existential types with explicit constraints in addition to inferred existential types. We keep track of existential types in cell `ex_types`, storing arbitrary constraints. For `LetVal`, we form existential types after solving the generated constraint, to have less intermediate variables in them. The first argument of the predicate variable  $\chi_K(\alpha, \gamma)$  provides an “escape route” for free variables, e.g. precondition variables used in postcondition. It is used for convenience in the formalism. In the implementation, after the constraints are solved, we expand it to pass each free variable as a separate parameter, to increase readability of exported OCaml code.

For simplicity, only toplevel definitions accept type and invariant annotations from the user. The constraints are modified according to the  $\llbracket \Gamma, \Sigma \vdash \text{ce} : \forall \bar{\alpha} [D]. \tau \rrbracket$  rule. Where `Letrec` uses a fresh variable  $\beta$ , `LetRecVal` incorporates the type from the annotation. The annotation is considered partial,  $D$  becomes part of the constraint generated for the recursive function but more constraints will be added if needed. The polymorphism of  $\forall \bar{\alpha}$  variables from the annotation is preserved since they are universally quantified in the generated constraint.

The constraints solver returns three components: the *residue*, which implies the constraint when the predicate variables are instantiated, and the solutions to unary and binary predicate variables. The residue and the predicate variable solutions are separated into *solved variables* part, which is a substitution, and remaining constraints (which are currently limited to linear inequalities). To get a predicate variable solution we look for the predicate variable identifier association and apply it to one or two type variable identifiers, which will instantiate the parameters of the predicate variable. We considered several ways to deal with multiple solutions:

1. report a failure to the user;
2. ask the user for decision;
3. perform backtracking search for the first solution that satisfies the subsequent program.

We use an enhanced variant of approach 1 as it is closest to traditional type inference workflow. Upon “multiple solutions” failure the user can add `assert` clauses (e.g. `assert false` stating that a program branch is impossible), and `test` clauses. The `test` clauses are boolean expressions with operational semantics of run-time tests: the test clauses are executed right after the definition is executed, and run-time error is reported when a clause returns `false`. The constraints from test clauses are included in the constraint for the toplevel definition, thus propagate more efficiently than backtracking would. The `assert` clauses are: `assert = type e1 e2` which translates as equality of types of `e1` and `e2`, `assert false` which translates as `CFalse`, and `assert e1 <= e2`, which translates as inequality  $n_1 \leq n_2$  assuming that `e1` has type `Num n1` and `e2` has type `Num n2`.

[TODO: implement the `assert` and `test` clauses!]

## Bibliography

- [1] Łukasz Stafiniak. A gadt system for invariant inference. Manuscript, 2012. Available at: <http://www.ii.uni.wroc.pl/~lukstafi/pubs/EGADTs.pdf>