

Parallel Type-checking with Saturating LVars

Peter Fogg Sam Tobin-Hochstadt Ryan R. Newton

Indiana University
{pfogg, samth, rnewton}@indiana.edu

Abstract

Given the sophistication of recent type systems, unification-based type-checking and inference can be a time-consuming phase of compilation—especially when union types are combined with subtyping. It is natural to consider improving performance through parallelism, but these algorithms are challenging to parallelize due to complicated control structure and difficulties representing data in a way that is both efficient and supports concurrency. We provide a solution to these problems based on the LVish approach to deterministic-by-default parallel programming. We extend LVish with a novel class of concurrent data structures: *Saturating LVars*, which are the first LVars to safely release memory during the object’s lifetime. Our new design allows us to achieve a parallel speedup on worst-case (exponential) inputs of traditional Hindley-Milner inference, and on the *Typed-Racket type-checking algorithm*, which yields up an $8.46\times$ parallel speedup on type-checking examples drawn from the Racket repository.

1. Introduction

Recent programming language advances often rely on sophisticated type systems [3, 4, 17, 31, 36], many of which incur a substantial computational expense at type-inference or type-checking time. In some cases, such as *Liquid Haskell’s* refinement types, it is possible to offload this work to an optimized (SMT) solver [17]. In other cases—occurrence typing, dependent typing, and gradual typing [35]—using an external solver is infeasible. Gradually-typed languages, for example, may employ uncommonly expressive type systems to capture idioms from dynamically typed programming. For example, Typed Racket combines subtyping with flexible union types, and the $(+)$ operation ends up with a type that combines several hundred distinct function signatures. As a consequence, the Typed Racket repository¹ contains individual files which take over five minutes to typecheck.

In the era of ubiquitous parallel hardware, one idea is to *parallelize* these computationally expensive phases to reduce the compile-edit-debug latency and enhance the software development experience. Yet there has been little work on parallelizing compilation of code below the granularity of a file or module, with the

¹<http://github.com/plt/racket>

exception of register allocation [38] and flow analyses [26]. Further, to the best of our knowledge, no prior work has parallelized type-checking algorithms specifically.

Most type checkers involve a unification process that contains latent parallelism but exhibits poor locality. A simple example is Hindley-Milner type inference, in which distinct expressions might be processed in parallel, but where each individual *type variable* can gain information from distant parts of the program (and therefore from different threads). Indeed, even in functional-language implementations of type inference (such as the Haskell-based implementation inside GHC), mutation is often used for constraining type variables, complicating parallelization further. Our goal is to parallelize despite these constraints, using linguistic abstraction to enforce disciplined, monotonic use of mutable state and minimize or eliminate unintended nondeterminism.

In this work we perform two experiments in parallel typechecking, and we use Haskell as our implementation language for writing parallel checkers. Because type checking is a constraint satisfaction problem, one strategy is to begin with a general constraint solver, extending it to deal with type checking. Yet this has not been done before, nor is it obviously easier than the approach we take here, which instead takes as its starting point standard implementation techniques for sequential type-checkers. Unfortunately, the mechanisms used to mutate type variables in these sequential checkers (State and ST monads) do not generalize safely to the parallel case while retaining referential transparency (which implies deterministic parallelism). Moreover, as described in §5.1, using immutable data structures with purely functional task parallelism (futures) is not a good fit for this problem.

LVars for Type Variables? Fortunately, there are a class of *synchronization variables* that are safe to share between computations in a functional language without compromising determinism. Single assignment variables, or *IVars* [2], are one early example in this class, supported in Haskell via the *Par monad* [23]. More recently, Kuper and Newton [18] introduced a generalization that enables deterministic, functional programs to synchronize on arbitrary *monotonic data structures*, called *LVars*. LVars enable a general form of deterministic-by-default parallel programming, implemented in Haskell by the “LVish” library². Previous work on LVars introduced general-purpose LVar data structures including: (1) lock-free collections with concurrent insertion (but not deletion), and (2) counters that increase monotonically. But application-specific LVars can be constructed as well; in particular, LVars would *seem* to provide a promising way to deal with type variables shared between threads. A custom LVar could capture the partial order implied by type unification. However, two problems arise:

1. All published examples of LVars respond to conflicting information by throwing an exception, which cannot be caught except in the `IO` monad.

[Copyright notice will appear here once ‘preprint’ option is removed.]

²<http://hackage.haskell.org/package/lvish>

- While LVars are a good fit for *And-parallelism*—where threads join information concurrently—they do not help with the *Or-parallelism* found in some type systems, where speculative, alternative additions of information must be considered.

Contributions In this paper, we solve these two problems and demonstrate the **first wall-clock parallel speedup on type inference**. Specifically:

- We introduce Saturating LVars (§4), adding the capability for both trapped-failure and memory reclamation—addressing a major limitation of previous LVar designs. We use Saturating LVars in the process of speeding up an implementation of Hindley-Milner type inference on some inputs.
- We develop a modular formulation of Or-parallel constraint systems, parameterized by an algebra for manipulating streams of partial solutions. We provide an efficient implementation of this algebra using *generators*. We explain this system in the simplified context of satisfiability problems (§5).
- We then scale this architecture to the type system of a full blown language (Typed Racket, §6) with a modestly widespread user base, achieving both speedups due to deforestation and parallel speedups.

2. LVars & LVish: Background

LVars generalize the earlier IVar model by allowing multiple writes. Where IVars simply signal an error upon writing to an already-full location, LVars allow the states to be *joined* in a monotonically increasing fashion according to a partial order on the possible states of the data structure. The state space (hereafter *lattice*) contains two distinguished elements \perp and \top —representing uninitialized and error respectively—along with a partial ordering \sqsubseteq . One way to increase the state of an LVar is through a *put* operation that takes the *least upper bound* of its current state and the argument to the *put*.

LVars also allow a restricted form of read via the *get* operation. Generalizing the blocking reads of IVars, this operation will block until the LVar’s state has reached one of a designated subset of the lattice’s elements, known as the *threshold set*. This set is, semantically, an implicit argument to *get*; it allows us only to observe that the LVar is *above* some element of the threshold set, rather than its precise state. The threshold set Q is required to be *incompatible*, that $\forall a, b \in Q, a \sqcup b = \top$.

As a simple example, consider the lattice of natural numbers ordered by the relation \leq . In the following program, two threads race to write to an LVar lv :

```
do lv ← newMaxIntLVar
  fork (put lv 1); put lv 2
```

Regardless of the order in which the threads write to lv , the join operation ensures that the final state of lv is “2”—the lub of both writes. As a result, we can freely share LVars between threads, safe in the knowledge that we will deterministically receive a result (or an error, in the case of the \top state), because *put* operations always *commute*.

In addition to thresholded *get* operations, changes to an LVar can be observed through *handlers* (callbacks). When we attach a handler to an LVar, it is called upon *each* change, and receives the new state (or, in some cases the delta itself). For the `maxIntLVar` shown above, we can attach a handler that is called each time the maximum is increased. With a container LVar on the other hand—such as a *Set*—the handler would be called on every element added to the *Set*, e.g.:

```
addHandler setLV ( $\lambda x \rightarrow \dots$ )
```

One interesting aspect of handlers is that `addHandler` must *commute with puts*. That is, upon adding the handler it fires for all past and future elements of the set above. Later in this paper we will see how the concept of handlers interacts with our proposed extension, saturating LVars.

Finally, LVars also offer the option of reading their full contents *exactly*, after they have been *frozen*. The `freeze` operation disallows further modifications, raising an exception if this occurs. For full determinism, freezing may only occur after a global barrier to avoid races between `put` and `freeze`.

LVars, in practice In the implementation of the LVish library, parallel computations are exposed through a `Par` monad. These computations have type `Par e s a`, where: a is the return value of the monadic `Par` computation; the s parameter keeps LVars from being shared between different parallel regions (like the `ST` monad); and the e type parameter documents the *effect signature* of the computation. For example, a function over `Ints` that executes in the LVish monad and also may *put* to an LVar, has this type:

```
foo :: (HasPut e) => Int -> Par e s Int
```

The `Par` monad is further equipped with various functions to launch parallel computations and extract their result:

```
runPar      :: Det e => ( $\forall s. Par e s a$ ) -> a
runParNonDet :: ( $\forall s. Par e s a$ ) -> IO a
```

These ensure that only deterministic (`Det`) combinations of effects are used from inside purely functional code, whereas nondeterministic combinations require `IO`. Both of these `runPar` variants are used to return *pure* Haskell results of type a . If one wishes to *return* LVars, they use `runParThenFreeze`, which uses the implicit barrier at the end of a `runPar` parallel region to guarantee a race-free freeze of the result:

```
runParThenFreeze :: (Det e, DeepFrz a)
                  => Par e NonFrzn a -> FrzType a
```

Freezing has no runtime cost. Rather, `FrzType` is a type-level function (type family) that “casts” the monotonic/mutable version of the LVar to a pure/immutable sister type. `FrzType` is associated with the `DeepFrz` class and implemented by each LVar in the library. (`NonFrzn` is a safety detail—placed in the s parameter to prevent mutating LVars that were frozen in other `runPar` sessions.)

A note on notation Effect signatures are important to the LVish library; some effects are only *conditionally* deterministic; for instance, canceling read-only futures is fine, but canceling a call to `foo` above would introduce an observable data-race. Yet effect signatures are not central to the way we use LVars in this paper, so we will abbreviate `Par e s Int` as `Par Int`, and mention effect constraints in the prose only where they are relevant.

3. And-Parallelism: Hindley-Milner Typing

We now have what we need to parallelize a type checking algorithm. We begin with what is perhaps the most well-known unification-based type inference algorithm: the Hindley-Milner system [8].

Sequential Hindley-Milner The Hindley-Milner algorithm operates by walking over an expression, generating constraints over type variables. These constraints are unified together to produce a final typing judgement for the term. An implementation that uses only immutable data would keep a store mapping type variables to types: `Map Var Type`. Recursive calls in the unifier produce partial maps that are *joined* together. This process is widely regarded as inefficient, and in practice even type checkers written in Haskell use a *mutable* representation of type variables. In this case, an explicit type variable store is unnecessary and monomorphic types can be defined as:

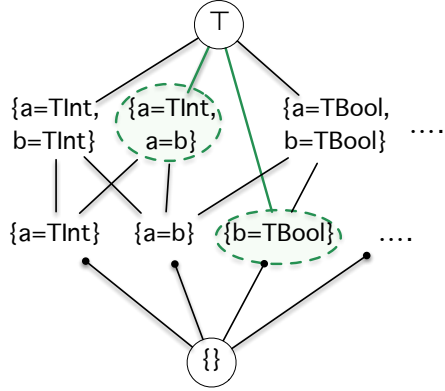


Figure 1. The partial order for the store containing all type variables used in a type-inference execution. The two highlighted nodes are *incompatible*—their lub is \top .

```
data Mono s = TVar Name (STRef s (Maybe (Mono s)))
  | TInt
  | TFun (Mono s) (Mono s)
```

Here the type variable is represented directly by its pointer to the mutable location. The `STRef` [21] allows a type variable to be imperatively updated in the `unify` function, as shown below.

```
unify :: Mono s → Mono s → ST s ()
unify t1 t2 = do
  case (t1, t2) of
    (TVar v ref, t) → do
      when (not (occurs v t)) (writeSTRef ref (Just t))
    (t, TVar v ref) → do
      when (not (occurs v t)) (writeSTRef ref (Just t))
    (TFun t1 t2, TFun t1' t2') → do
      unify t1 t1'
      unify t2 t2'
    (TInt, TInt) → return ()
    _ → error "can't unify"
```

The `infer` function simply walks over the expression, unifying type variables as they are found. For brevity, only the application case is shown.

```
infer :: Env s → Term → ST s (Mono s)
infer env expr = case expr of
  ...
  App e1 e2 → do
    fnT ← infer env e1
    arg ← infer env e2
    res ← freshTVar ()
    unify fnT (TFun arg res)
    return res
```

Alas, in exchange for increased performance in the single threaded case, it would appear that `STRefs` have destroyed our opportunity for deterministic parallelism in this implementation!

Exploiting parallelism Fortunately, we can parallelize the algorithm by switching from `STRefs` to `LVars`, and can even asynchronously share type constraints *between threads*. The collection of type-variable constraints accumulated during type-checking forms a partial order under unification (Figure 1). We can model this state either with a single `Map LVar` or with an `LVar` per type variable (`TyVar`)—we use the latter approach here.

A `TyVar` is an application-specific `LVar` that is either empty (unconstrained) or filled with a type, and when additional types, containing type variables, are joined into the `TyVar` via `put`, unification is invoked recursively as a callback. With `TyVars`, the previous definition of `Mono` changes to include:

```
data Mono s = TVar (TyVar s) | ...
```

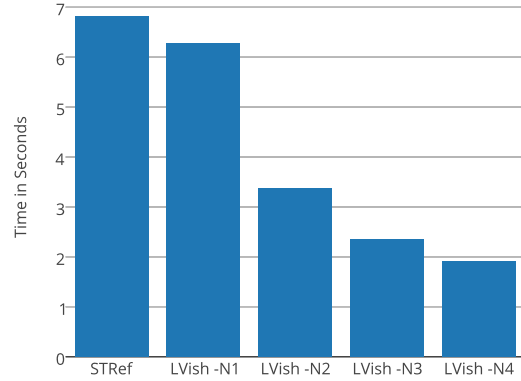


Figure 2. Hindley-Milner type inference on a generated term with roughly a million nodes. The sequential (`STRef`) version is compared against the `LVish` version. Note that even on one core, the `LVish` version is still *concurrent*, sharing constraint information between concurrent computations via `LVars`. Experiments were performed on a desktop-class Intel Xeon i5-3470, with GHC 7.8.3 and `+RTS -qa`.

Because `TyVars` are simply a mutable pointer to a piece of immutable data, the `TyVar` abstract datatype is easy to implement—in contrast with container `LVars` that require complicated lock-free algorithms.

Much of the type-checking code is fundamentally unchanged, beyond replacing `writeSTRef` with the analogous operation on `TyVars` and switching type signatures for their `Par` version. After that refactoring, we can inject parallelism in the application case of the `infer`, to then infer the types of both subexpressions in parallel:

```
infer :: Env → Term → Par Mono
infer env expr = case expr of
  App e1 e2 → do
    (fnT, arg) ← par2 (infer env e1) (infer env e2)
    res ← freshTVar
    unify fnT (TFun arg res)
    return res
```

Here the `par2` combinator simply executes two actions, returning two values (a `fork/join` construct). We can also unify function types in parallel, which is a `fork` with no `join`:

```
unify :: Mono → Mono → Par ()
unify (TFun t1 t2) (TFun t1' t2') = do
  fork (unify t1 t1')
  unify t2 t2'
```

Note that unification called for (monotonic) side effects only. This approach yields a speedup when used to implement a micro-ML calculus and run on a synthetic benchmark (see Figure 2). (For type-checking benchmarks drawn from actual code, see §6.) In this case, the benchmark is a large program with 1000 copies of a known-exponential nested-let expression. This is only a small proof of concept—for real, large programs with Hindley Milner inference, the challenge will be finding problems that take long enough in practice (relative to the amount of memory they read) that they are worth parallelizing. Nevertheless, with algorithms such as Hindley-Milner—which perform well in the average case but have dismal worst-case performance—we believe it is worth researching parallelism as an “insurance policy” to ameliorate these worst case outcomes.

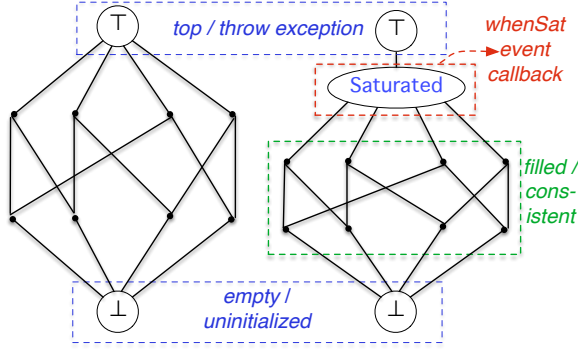


Figure 3. Any valid LVar lattices is turned into a Saturating LVar by adding an extra, penultimate state.

4. Saturating LVars: Trapped Failure

There is one problem with the formulation of Hindley-Milner type inference in the previous section—what if two types do *not* unify? If type-checking fails we would like to simply return `False` or `Nothing`. But with implementation in the previous section this failure instead appears as incompatible puts, e.g. putting `TInt` and `TBool` to an LVar. Further, in LVar-based programs, adding contradictory information to an LVar always triggers the \top state, which in previous implementations of LVish meant throwing a Haskell exception.

What is wrong with throwing an exception? Answering this requires a bit of background. Haskell enables purely functional but *partial* programming, and Haskell’s exception semantics [24] requires that exceptions be handled only in the `IO` monad to retain referential transparency. In this case, it is important that we keep the type checking phase *out of* the `IO` monad. Because we aim for a deterministic type checker, we should either use only determinism-safe features (not `IO`) or reduce the amount of “trusted code” that uses unsafe features that may introduce nondeterminism—most especially avoiding Haskell’s infamous “`unsafePerformIO`”.

In fact, there are significant benefits to strictly deterministic compilers, generally. The Nix package manager and NixOS operating system [9] have demonstrated the benefits that accrue from compilation being a mathematical function from bits to bits. Writing a compiler and type-checker in Haskell with no `IO` (outside of reading and writing files) is one way to achieve this goal.

Keeping failures in Par In order to avoid the exception handling problem, we must capture and respond to type checking failures *within* the `Par` monad. That is, when type variables gain conflicting information, we want to simply return a value indicating no valid substitution exists.

To enable trapped failures, we introduce *Saturating LVars*, defined as an LVar whose lattice structure includes an additional state `Sat`. Following the semantics for LVars in POPL’14 [20], such an LVar is given by a five-tuple $(D, \sqsubseteq, \perp, \top, Sat)$, extended to include the designated saturation state as well as the usual set of states (D) , partial order, and designated bottom and top states. It should further hold that:

$$\perp, \top, Sat \in D, \perp \neq \top, Sat \neq \top$$

$$\forall d \in D, (\perp \sqsubseteq d \sqsubseteq \top) \wedge (d \sqsubseteq Sat \vee d = \top)$$

An example lattice extended with the `Sat` state is pictured in Figure 3. While this convention is simple, its ramifications are not:

1. The only usable (incompatible) threshold set for performing blocking read or adding a handler³ is $\{Sat\}$.
2. A saturating LVar’s state moves monotonically up the lattice, but it does *not* monotonically gain information (bits). Notably an LVar in the saturated state can be represented by as little as one bit. This means that saturating LVars are the first LVars **that can release memory** during their lifetime.
3. Computations whose *only* effect is to write to a Saturating LVar can be *cancelled* if that LVar saturates (§4.1).

Because of the first restriction, saturating LVars become effectively *write only*. Further, all Saturating LVars can provide the following operations:

```
class DeepFrz lv => SatLVar lv where
  saturate :: lv -> Par () -- Force it to Sat state
  whenSat  :: lv -> Par () -> Par ()
  isSat    :: FrzType lv -> Bool
```

This interface provides the ability to force LVars to saturate, respond to saturation, and test for saturation after a parallel computation is complete.⁴

Of course, every saturating LVar provides specific methods outside of the common interface—methods for constraining a type variable, inserting into a map or set, and so on. But even the common API is enough to consider some use cases.

Example use-cases Consider an application-specific LVar representing constraints upon a single variable. The `Sat` state would correspond to conflicting constraints (i.e. *failure*). If we collect those individual LVars into a container, such as a `Map` LVar, it could represent a complete environment of variable assignments. Thus we would expect that the *entire environment fails* when any entry does. Indeed, this is possible with the API we described above—before inserting each variable in the `Map`, we attach a `whenSat` handler which in turn calls `saturate` on the entire environment, propagating the failure, as in the following:

```
do env <- newEmptyMap
  ...
  cv <- newConstrainedVar
  whenSat cv (saturate env)
  insert vr cv env
```

In fact, there’s more that a library can do to enable efficient compositions of LVars and saturating LVars. For example, we have produced a new, modified LVish library that supports saturating LVars and in this library we provide:

- a `SatMap` data structure, which is a *single* LVar that maps keys onto pure Haskell values that are instances of `PartialJoinSemiLattice`. That is, multiple puts are allowed on the same key, and are joined, but the `join` function may fail, saturating the entire `SatMap`.
- a `FilterSet` data structure that takes advantage of saturated LVars in a different way—it represents a dynamic collection of not-yet-failed saturated LVars. We observe that the type `Set (Maybe a)` is isomorphic to `(Bool, Set a)`, that is, there’s no need to store *each* element which has saturated. If we are interested in collecting `SatLVar`’s that have *not* failed, then failed LVars can be discarded at runtime, freeing memory and shrinking the set’s physical size.

³ Actually, there is a safe way to add handlers that are notified of *every* put to the `SatLVar`, but it requires that the handlers be attached *at the point the LVar is created*, which prevents `addHandler` racing with `saturate`.

⁴ The LVar can only be tested for saturation with `isSat` *after* a parallel region has ended and it is in a “frozen” state. Freezing LVars is covered in detail in [20].

insert(10)	OLS Regression	R^2 goodness-of-fit
Map LVar	19.5ns	0.991
SatMap	18.4ns	0.993
Set LVar	18.5ns	0.989
FiltSet	91.1ns	0.990

Table 1. Microbenchmark: the cost of creating a new structure and inserting ten new `Int` elements. The cost of this (comparatively) cheap operation is measured by varying the number of iterations of benchmark, and computing a linear regression between iterations and cycles (above). All measurements are from the *desktop* platform described in §6.3.

insert/sat/insert	Set LVar	FiltSet
cycles	17838	15160
bytes alloc	14733	14128
bytes copied	759	115

Table 2. Microbenchmark: the cost of inserting 10 `Counter` elements in a set, saturating the previous 10, and repeating N times. The `LVar.Set` version must store the data as a set of nested LVars to enable saturation of the inner variables. The `FiltSet` directly supports multiple assignments to a key, so requires one LVar rather than $10N + 1$. Above we regress N against cycles, and below we regress against bytes allocated and bytes copied during garbage collection. This verifies that the while the `FiltSet` benchmark allocates $O(N)$ memory, it releases memory as it goes.

As we will see, we can use `FiltSet` to accumulate results from a search process when dealing with type systems that include disjunction. Or, as a simpler example, consider implementing a program analogous to the following query using a data structure of type `(FiltSet SatCounter)`:

```
SELECT MEDIAN(SIZE) FROM CLASSES WHERE SIZE < 30
```

While traversing a list of students in parallel, we can maintain a set of counters that are updated with `fetch-and-add`, and are set to saturate upon hitting a tally of 30. Over-threshold counters would automatically be removed from the set, leaving only those which are under the threshold at the end of the `runParThenFreeze` call.

In the Hindley-Milner type-checking case, we use saturating LVars for each individual `TyVar`. Saturation does not come into play in typing *well-typed* terms as in Figure 2. But when discovering that a term is ill-typed, saturation (and as we will see below, *cancellation*) are relevant. Further, in other type systems, the `SatMap` structure is useful for representing environments containing constraints, and a `FiltSet` can serve as the accumulator when searching for a valid environment. These two data structures are part of the parallel-type-checking toolkit we provide in our new library; the microbenchmarks in Tables 1 and 2 show their performance relative to more basic LVar counterpart data structures.

4.1 Saturation and Cancellation: a safe idiom

Because saturating LVars are write-only in parallel regions, a natural question arises: will threads continue blindly adding information to an LVar that has *already* saturated? If the continuation of the computation that saturates the LVar performs no other observable side effects, it is better to *cancel* remaining computations in response to saturation. But is this safe to do?

Cancellation is a feature supported by LVish and explored in previous work [19]. In that work, however, cancellation was only safe for LVish computations with `get` but never any kind of write effect: e.g., `put`, `freeze`, or `saturate`. That restriction rules out canceling upon saturation, because to have saturation events in the first place, we must perform *writes* to an LVar! The insight here is that, based on the lattice structure of saturating LVars (Figure 3), after

saturation, all *further* writes have no observable effect (including not triggering handlers). Indeed, we can phrase this assertion in the terms of the previously published LVar semantics [20]:

Conjecture 1 (Safe-Cancellation). *Every LVish program with store S , such that for all writable LVars l , $[l \rightarrow \text{Sat}] \in S$, is observationally equivalent to the program `return ()` or \perp (divergence).*

If the author of a particular application is amenable to converting potential nonterminating outcomes into cancelled ones, then they could elect to leverage this property to cancel useless computations. We restrict this the case of writing to a single saturating LVar for simplicity. But how then do we enforce that a given sub-computation can *only* write that LVar, while reading from arbitrarily many other LVars? To accomplish this, we can use LVish’s effect-tracking capabilities to formulate a safe `withSatLVar` idiom, which runs a block of code with a `ReadOnly` effect signature, but with an “escape hatch” for modifying *only* the designated saturating LVar:

```
withSatLVar lv (λ modIt →
  do (k,v) ← ...readOnly...
    modIt (insert k v) -- Send out a write to lv
    ...readOnly...)
```

In this example, `lv` is a `SatMap`, and the `modIt` function allows executing code in a separate effect environment that allows writes, but which can operate *only* on `lv` and no other LVar. Like most programs with tracked effects, the type of `withSatLVar` is more complicated than the code that uses it. The type `withSatLVar` makes heavy use of the “e s” parameters which we have been eliding when writing `Par a` rather than `Par e s a`. Here we show that full type:

```
withSatLVar :: (SatLVar lv, ReadOnly e1, Det e0)
  ⇒ lv s1
  → ((∀ s0. lv s0 → Par e0 s0 ())
     → Par e1 s1 a)
  → Par e2 s1 (Maybe a)
```

Here the higher-rank type ensures that the computation passed to `modIt` really can touch *only one* LVar. In exchange for this harsh restriction, the effect environment `e0` need not be in any way connected to `e1`—but it must remain deterministic. Note that because of the cancellation possibility, a `Maybe` value is returned. It is `Nothing` if `lv` saturates and the computation is cancelled. In the case of Hindley-Milner type inference, `Nothing` corresponds to a type error found by the type checking algorithm.

5. Or-Parallelism: Satisfiability solvers

With LVish plus the saturating LVar extension, we’ve acquired the first tools in our parallel type-checking toolbox, enabling us to handle parallel *conjunctions* over constraint-generating computations. Indeed, Hindley-Milner type inference required only conjunction, never disjunction. In this section we begin to address a broader—and more expensive in practice—class of type checking algorithms: those with *Or-parallelism*. This is challenging, because we cannot directly employ LVars the way we did in the previous section (mapping each type variable to exactly one LVar and updating it destructively). Nevertheless, Or-parallelism is a core feature of Typed Racket, the type system that is our main target in this paper (§6). Yet rather than dive directly into Typed Racket in this section, we first introduce the implementation techniques using a smaller example with conjunctions and disjunctions: satisfiability (SAT).

Parallel Constraint Solving Type systems are a particular flavor of constraint problem. Indeed, if we view parallel type checking as a parallel constraint satisfaction problem, we can look for guidance from previous work on parallel logic programming [6, 13] and parallel constraint solvers [11, 27, 37]. Unfortunately, the data-structures and synchronization strategies employed in these works

```

generic :: ∀ a b. Eq a ⇒
  Ringlike (Vector (Var,a)) b → Tree a → b
generic Ringlike{mkNum,mkZer,add,mul} tr0 =
  case tr0 of Left at → loop1 at
            Right ot → loop2 ot
where
loop1 :: AndTree a → b
loop1 (Leaves vec) = mkNum vec
loop1 (And vec) | null vec = mkNum empty
                | otherwise = mul (map loop2 vec)
loop2 :: OrTree a → b
loop2 (Or vec) | null vec = mkZer
               | otherwise = add (map loop1 vec)

```

Figure 4. The generic solution for satisfiability, parameterized by a `Ringlike` object.

are extremely specialized to the constraint system being solved: for example, SAT solvers have developed a large body of specialized data structures and parallelization strategies [14, 15].

We are not here concerned with specialized search strategies—e.g., our goal in this section is not to implement an efficient SAT solver. Rather we study data-structure trade-offs and parallel control-flow trade-offs when exhaustively searching a space of possibilities using techniques that apply to *any* constraint domain that can be formulated as an LVar. We will then apply the same techniques to Typed Racket in §6.

In the case of satisfiability, we have a constraint domain that consists of simple variable assignments, e.g., $x=4$, closed under conjunction and disjunction. Thus the input to our algorithm is a term such as:

$$((x = 3, y = 4) \vee (y = 3)) \wedge (y = 3, z = 9)$$

Represented in our Haskell implementation by the following data structure, which enforces a normal form where `And` and `Or` alternate, but both are N -ary rather than binary:

```

type Tree a = Either (AndTree a) (OrTree a)
data AndTree a = And (Vector (OrTree a))
                | Leaves (Vector (Var,a))
data OrTree a = Or (Vector (AndTree a))

```

A simple, compositional solution must represent a *stream of substitutions*, each representing possible bindings for variables within the (sub)term. These solution streams can be combined by concatenation (`Or`) or by joining pairs of solutions drawn from the cartesian product of two streams (`And`). In fact, the streams form a ring-like structure, and we can formulate a simple generic solution abstracted over a set of methods matching the following signature:

```

data Ringlike leaf -- "leaves" of our computation
                  elem -- elements of our ring
= Ringlike
{ mkNum :: leaf → elem
, mkZer :: elem
, add :: Vector elem → elem
, mul :: Vector elem → elem
}

```

This provides a simple algebra for solution streams. The generic code that walks the `Tree` and calls the `Ringlike` methods is listed in Figure 4, and is used by all the implementations we discuss below.

5.1 The Simplest Stream Algebra

In a sequential Haskell implementation, the natural representation of solution streams is as a lazy list of variable assignments (`Maps`), parameterized over the type of values variables range over, “a”:

```

type Env a = Map Var a
type Sol1 a = [Env a]

```

And the algebra over these streams is implemented by:

```

listStrms :: Eq a ⇒ Ringlike (Vector (Var,a)) (Sol1 a)
listStrms = Ringlike
{ mkNum = λvec → maybeToList (foldlM insert empty vec)
, mkZer = []
, mul = foldl1 (λs1 s2 →
                catMaybes [ joinEnvs env1 env2
                          | env1 ← s1, env2 ← s2 ])
, add = foldl' (++) []
}

```

Here the type argument “`Vector (Var,a)`” represents a block of variable assignments, and `insert` is a function that gives the block an interpretation in terms of `Sol1s` (inserting the variables into the `Map`, and failing if there is a conflict using the `Maybe` monad). The cartesian product operation above creates a list of many `join` computations, which could be evaluated in parallel. In fact, we attempted to parallelize in the standard Haskell way by adding `parList` or `parBuffer` annotations to this list, either before or after the `catMaybes` call. Unfortunately, this does not yield a parallel speedup (either for satisfiability or full Typed Racket), because the genuine parallel work is too entangled with book-keeping on lazy lists.

5.2 A Parallel Stream Algebra with Generators

As we will see in §6, list-based streams incur a lot of overhead—intermediate lists are assembled and deconstructed repeatedly. Further, the aggressive fusion optimizations performed by GHC and its libraries **cannot eliminate operations like cartesian product**.

Fortunately, there are more efficient ways to represent streams, in particular as *generators*. Generators have a long history as a control mechanism in programming languages⁵. A generator takes a partial answer and a continuation; it modifies, tests, or bifurcates the partial answer; and then passes one or more answers on to the continuation.

```

type Cont = PartialAns → [PartialAns]
type Generator = Cont → Cont
              = Cont → PartialAns → [PartialAns]

```

Generators can be composed without creating intermediate lists—only the final step allocates `[PartialAns]`. Indeed, generators, formulated in terms of continuations, have been used for this deforestation benefit in many contexts. Yet there has been little work on their use in parallel programming⁶. We can, for example, define our answer type to be a computation in the LVarish `Par` monad, which gives us the following solution type for satisfiability problems:

```

type Cont a = Env a → Par ()
type Sol3 a = Cont a → Cont a

```

```

contBased :: (...) ⇒ Ringlike (Vector (Var,a)) (Sol3 a)
contBased = Ringlike { ... }

```

A solution stream—the element type of our `Ringlike`—becomes a function of the form $(\lambda k w \rightarrow \text{action})$, where `action` constrains the variable assignment, `w`, in one or more ways and passes each variant on to the continuation `k`.

Further, note that each computation, `Par ()`, does *not* return a value. As in CPS (continuation passing style), when using this formulation we extract a result by attaching a final continuation that inserts into an output `Set` or `FiltSet` LVar.

With this definition for `Sol3`, we can see that the *zero* for the algebra is $(\lambda k w \rightarrow \text{return } ())$, which drops the partial assignment

⁵ First appearing in the language `Alphard` in the mid 1970s [29], and later in `CLU` [22] and `Icon` [12]. A clear explanation of generators can be found in [1].

⁶ One example in the literature is the related concept of *push arrays* [5] used in data-parallel programming.

w on the floor, *not* calling the continuation. Likewise the *one* value is $(\lambda k w \rightarrow k w)$. The disjunction, or add operation is the obvious place to add parallelism:

```
add s t = \ k w \rightarrow
  do fork (s k w)
     t k w
```

This parallel-OR duplicates the continuation, passing the incomplete answer to alternative code paths that extend it in different ways. Here we show the binary version of the operator, but the N -ary version in our implementation is a straightforward generalization which exposes the parallelism as a parallel loop rather than a single `fork`.) Further, we thus far assumed immutable environments, such that `w` does not need to be *copied* before being sent to different destinations (`s` and `t`). We will change that in a moment.

AndPar: first technique The sequential version of binary conjunction with immutable environments is:

```
mul s t = \ k w \rightarrow s (\lambda x \rightarrow t k x) w
          = \ k \rightarrow s (t k) = s o t
```

Indeed, there is no obvious opportunity to parallelize this as long as environments are immutable. The continuation transformers are composed, but `w` is threaded through linearly. Likewise `mkNum`, which handles the Leaves of the `AndTree`:

```
mkNum vec = \ k w \rightarrow k (foldl constrain w vec)
```

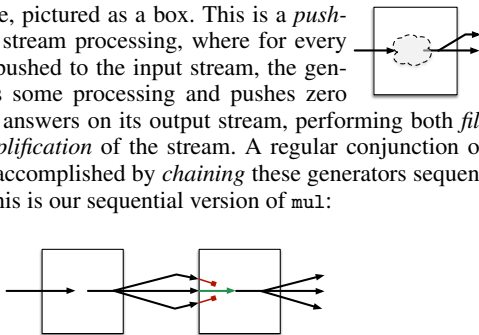
With the immutable definition of `Env` as `Map Var a`, we must fold the constraints into the environment *before* sending it along to the continuation. It is not possible to extract parallelism here at the level of the `Par` monad. (As in the previous section, it would be possible—but not profitable—to *spark* the `foldl` computation, attempting very fine-grained parallelism within the updates to a `Map`.)

If we use LVars to represent the environment, we retain the generator design but and-parallelism becomes more feasible. We change the `Env` to an LVar, such as `SatMap a`, and each generator modifies this environment *as an effect*. Generators representing conjunctions only perform these effects and *always* pass the same environment pointer on to their continuation that they receive. For example, $(k w)$ below:

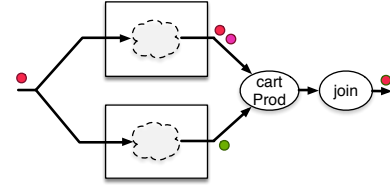
```
mkNum vec = \ k w \rightarrow do mapM_ (constrain_ w) vec; k w
```

Depending on the size of `vec` it is possible to `fork` the entire `(mapM_ ...)` expression or to turn it into a parallel loop.

AndPar: second technique There is another option for parallelizing and, but it may result in repeatedly traversing constraints that are already “settled”. We introduce it visually first and then in code. Let us visualize each continuation transformer as a processing stage, pictured as a box. This is a *push-driven* form of stream processing, where for every partial answer pushed to the input stream, the generator performs some processing and pushes zero or more partial answers on its output stream, performing both *filtration* and *amplification* of the stream. A regular conjunction of disjunctions is accomplished by *chaining* these generators sequentially. Indeed, this is our sequential version of `mul`:



The green arrow corresponds to a partial answer “passing the test” and moving on to the next round. Our second technique for performing And-parallelism is to replace the above picture with:



Here we take each partial answer, *duplicate* it, and feed it through both machines in parallel. The top and bottom boxes may or may not contain disjunctions. All partial answers that make it through the gauntlet on the top, are joined with all answers on the bottom, and, if the join succeeds, passed on. In code, we write this as:

```
mul s t = \ k w \rightarrow
  do s1 \leftarrow newEmptySet
     s2 \leftarrow newEmptySet
     s3 \leftarrow cartesianProd s1 s2
     forEach s3 (\ (a,b) \rightarrow case joinMaybe a b of
                               Nothing \rightarrow return ()
                               Just w2 \rightarrow k w2)
  fork (s w ('insert' s1))
     t w ('insert' s2)
  return ()
```

This uses LVars to accumulate the solutions from each branch, and to take their cartesian product (a monotonic operation, and a standard on container LVars). There is a delicate trade-off, however, in applying this technique: first, because the input answer ‘`w`’, is passed to both branches, all the joins we perform redundantly join the (obviously compatible) information in ‘`w`’ with itself. Second, we have now removed some of the deforestation benefit of generators by accumulating the partial answers in set LVars. Nevertheless, we in the next section we will see that this parallel-and technique is quite effective in typing some Typed Racket programs.

6. Typed Racket Type Checking

Typed Racket [35] is a typed version of Racket [10] that uses gradual typing [30, 34] to integrate with untyped Racket. In this paper, we consider only the type checking of typed programs.

Typed Racket’s type system includes a number of features which combine to make type inference difficult. First, Typed Racket supports subtyping, which is used widely in a variety of ways in Racket programs. Second, types include non-disjoint union types, so that $T <: (\cup S T)$. Third, overloading on function types is supported with ordered intersection on function types [32]. Fourth, Typed Racket supports arbitrary equi-recursive types, used for modeling even simple data structures such as lists.

As a result, Typed Racket, following many other recent languages, does not attempt complete whole-program type inference in the fashion of Standard ML. Instead, it employs so-called *local type inference* [25] along with bidirectional type checking.

In this setting, the central inference problem is choosing an instantiation of type variables when a polymorphic function is applied to concrete arguments. For example, when a Typed Racket programmer writes:

```
(map add1 (list 1 2 3 4))
```

we need to infer that `map` should be instantiated with the types `Integer` and `Integer`.

As a result, the type inference problem is somewhat simpler than in a global type inference setting. In the above case, if `map` has the type $\forall \alpha \beta. (\alpha \rightarrow \beta) \times \text{listof}(\alpha) \rightarrow \text{listof}(\beta)$, the inference algorithm must find a substitution for α and β that makes `Integer` \rightarrow `Integer` a subtype of $\alpha \rightarrow \beta$ and `listof(Integer)` a subtype of `listof`(α).

However, the inference problem in Typed Racket is made more complex than this simple example by several factors. First, type constraints in inference can involve subtyping, not just equality. Second, Typed Racket produces very large types in several circumstances—when providing extremely precise specification of function behavior [32] and when inferring a type for large blocks of constant data. As a result of these and other issues, Typed Racket is known to have slow typechecking. In particular, some pathological cases can have typechecking times measured in minutes. As mentioned in the introduction, one case which we use as a benchmark was removed from the Typed Racket test suite since it took too long to run.

6.1 The core algorithm

The core of the inference algorithm is an extended form of the Pierce and Turner [25] algorithm, which handles union types, recursive types, and function overloading. The fundamental idea is that we grow a set of constraints on the type variables to be inferred based on the actual types of the arguments provided—in the map above the actual arguments are `Integer → Integer` and `listof(Integer)`.

Each constraint is a pair of types: an upper and a lower bound. Two constraints can be combined by joining the lower bounds (represented by the \cup operation) and taking the meet of the upper bounds (meets cannot always be represented exactly in Typed Racket, requiring some approximation.) Inference fails if the constraints cross. To solve $S <: (\mu X.T)$, the algorithm must *unroll* recursive types; to ensure termination, the recursive solver must also keep a *seen* set, so that if, while unrolling, the same $S <: T$ is encountered again, the algorithm terminates successfully.

The additional complexity, and source of or-parallelism, comes from handling union types and overloaded function types. To make a type A a subtype of $(\cup B C)$, it must merely be a subtype of *one* of B or C . Therefore, the standard sequential algorithms as currently implemented in Typed Racket simply tries to solve $A <: B$, and if that fails, tries $A <: C$. Or-parallelism then enters by trying both of these possibilities simultaneously, succeeding if one succeeds. For function overloading, which is modeled as intersection in Typed Racket, we simply consider the dual problem, with similar implications for parallelism.

Inference in Typed Racket also has the possibility for and-parallelism. If we wish to constraint (A, B) to be a subtype of (C, D) , this implies a pair of constraints, both of which must succeed for a full solution.

6.2 Implementing Typed Racket Inference

We first performed a direct port of the `infer` implementation in the Typed Racket source code. This function takes the names of type-variables to constrain, as well as the relevant types for type-checking a polymorphic function invocation:

```
infer tvars actualTys formalTys resultTy expectedTy
```

In order to perform parallelization experiments using LVish, we port the code and the grammar of types to Haskell with one omission of functionality—we omit variable arity functions, which do not occur in our benchmarks. Also, while we keep the structure of the code the same in this conversion, we substitute some data structures with idiomatic Haskell counterparts (replacing lists with sets or maps in places). We refer to this as the “Pure/Seq” version of the program—purely functional, non-monadic, and sequential.

Refactoring for parallelism Next we rewrite the algorithm in monadic style and abstract the core constraint-gen (`cg`) recursion so that it returns a solution stream as in §5, and factors out the corresponding methods for conjunction and disjunction. Thus it is possible to use the same core algorithm with different *evaluation*

strategies, including sequential or parallel versions. Each implementation of the type checking solution algebra provides the following functions, which we have given names more appropriate to the task at hand:

- `goodsofar` (one) – result indicating no conflicts observed at this point in the search
- `constrain` (`mkNum`) – add an upper and lower bound constraint on a type variable, apply this to all partial solutions processed
- `blowup` (zero) – result indicating a conflict found
- `orSplit` (`add`) – test N alternative ($S <: T_i$) subtyping constraints, where $i \in [0, N)$
- `andPar` (`mul`) – join constraints with (optional) parallelism

Then, using the above, the following is a subset of the cases in the heart of the subtype checking algorithm, showing each possible behavior:

```
-- Make s a subtype of t:
case (s,t) of
(s, t) | (s, t) ∈ seen → goodsofar
(s, t) | s == t       → goodsofar
(s, t) | subtype s t  → goodsofar
(_, Top)              → goodsofar
(Var x, t) | x ∈ xs → constrain bot x (demote vs t)
(s, Rec _ _)         → cg s (unfold t)
-- N-way Or-parallelism:
(s, Union ts)       → orSplit cg s (elems ts)

(Pair a b, Pair a' b') → andPar (cg a a') (cg b b')
...
_ → blowup -- s cannot be a subtype of t
```

Solution strategies We implement this parallel algebra of solution streams while leaving knobs to toggle and-parallelism and or-parallelism independently at compile time. The `andPar` function precisely resembles the code for `mul` in §5.2. Or-parallelism becomes a standard parallel `forEach` from the LVish library:

```
orSplit :: (a → a → Solution) → a → [a] → Solution
orSplit msg doConstraints s ts = λ k varmap →
  parForEach ts (λt → doConstraints s t k varmap)
```

Note that the `parForEach` is an asynchronous operation—it forks work and returns immediately, without any join.

Testing While we do not run directly on Typed Racket source programs, we *traced* calls to Typed Racket subtype checking procedure, generating a log of over 50,000 test cases that we validate our Haskell implementations (sequential and parallel) against, confirming all cases type check or fail as appropriate. Further, because we know the typing constraints should form a partial order, we randomly generate types with the QuickCheck library and test lattice properties—e.g., that everything is above bottom, or the lub of two constraints is above each of them.

6.3 Typed Racket Evaluation

For our implementation of the core of the Typed Racket inference algorithm, our evaluation focuses on two different demanding inference problems. First, we consider the case mentioned in the introduction—a small function that takes minutes to check. Second, we consider checking large constant data against a small type. These are both representative problems that have been identified as the most serious performance problems for Typed Racket.

“Bigcall”: Higher order functions over extremely polymorphic inputs Typed Racket supports both polymorphism and overloading, and when combined, these can produce computationally-intensive inference problems. The most significant of these is the

following⁷, designed as an example of Typed Racket’s *variable arity polymorphism* [33].

```
(: map-with-funcs
  (All (b a ...)
        ((a ... -> b) * -> (a ... -> (Listof b))))
(define (map-with-funcs . fs)
  (lambda as
    (map (lambda: ([f : (a ... a -> b)]) (apply f as))
         fs)))
((map-with-funcs + - * /) 1 2 3 4 5)
```

This function consumes a variable number of functions, bound to the list `fs` and then a variable number of arguments, bound to the list `as`. It then applies each function `f` from the list to *all* of the `as`. We then apply `map-with-funcs` to a few arithmetic functions and apply the result to numbers. The result is a list containing the sum, difference, product, and division of all five numbers (Racket’s numeric operations all support arbitrarily many arguments).

The sequence of arguments `fs` is described in Typed Racket using variable-arity polymorphism. Since we omit this portion of the algorithm, we instead consider versions of `map-with-funcs` that consume 1, 2, 3, or 4 arguments: i.e., `(map-with-funcs + - *)` is the three argument case. We name these *bigcall*(1) through *bigcall*(4) for brevity.

Solving this inference problem requires handling several type variables, each of which is jointly constrained by all the arguments, but more importantly, Typed Racket provides very large overloaded types to give precise specifications to numeric operations such as `(+)`, which has *hundreds* of possible types [32]. Since the type of each arithmetic operator is an intersection, any choice of a single overload for one can be combined with any choice of an overload for another input, resulting in a combinatorial explosion of possibilities. Thus type checking `bigcall`(4) takes many minutes to complete in Typed Racket.

“Treecall”: Dealing with large constant data The second challenge we consider is that of large constant data. Typed Racket supports flexible and precise types for structured data in s-expression format. If a large constant is present in a program and no extra annotation is provided, it will therefore infer the *most precise type*, which can be the same size as the data itself. When a polymorphic function such as `map` is applied to this data structure, inference must process this large type.

To simulate this in a controlled fashion, we designed a benchmark which is the equivalent of applying the following function to progressively larger inputs consisting of trees of symbols and strings.

```
(define-type (Tree A) (Rec X (U (Leaf A) (Pair X X))))
(: leftmost : (All (A) (Tree A) -> A))
(define (leftmost t)
  (if (pair? t) (leftmost (car t)) (leaf-val t)))
```

Hence we call this “treecall”, because it calls `leftmost` on a constant tree datum. `Treecall`(*N*) corresponds to applying `leftmost` to a tree of depth *N*, i.e. with 2^N leaves. In a language with a rich macro system like Racket’s, large compile-time data is a reality, and is currently a Typed Racket performance problem.

6.3.1 Benchmark Results

We evaluate `treecall` and `bigcall` on one desktop-class and one server-class system, with one Intel Xeon i5-3470, and two Xeon E5-2670 CPUs, respectively (4 and 16 cores). All experiments were run using GHC 7.8.3 and all data points are the median of 5 trials. We distinguish two different kinds of run, where we generate either *all* substitutions, or just the first, which we explain further below.

⁷Taken from Typed Racket’s online tests at <http://git.io/ve4PJw>

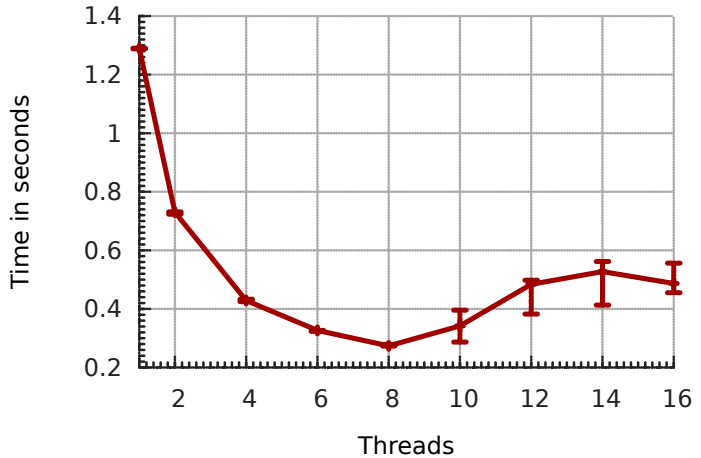


Figure 5. `Treecall`(16), `ParAnd` case: here scaling stops at eight cores. We plot time rather than parallel speedup and include minimum and maximum times across all trials in the error bars. In this way, we can see how runtimes become chaotic after scaling stops.

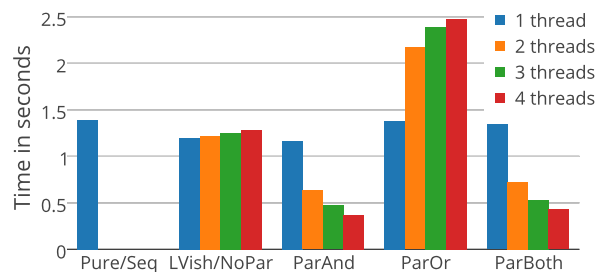


Figure 6. `Treecall`(16): Large-constant benchmark of size 2^{16} , desktop platform. Here Or-parallelism is not useful, because in each binary-Or, one branch is always an immediate dead-end.

Treecall results

Figures 6 and 5 show `treecall`(16)—the result of applying `leftmost` to a constant of size 2^{16} on the desktop and server platforms, respectively. In this benchmark there is only one solution, so “all vs. first” solution is immaterial. Because of the union in the list type, `(Rec X (U (Leaf A) (Pair X X)))`, there are many *apparent* opportunities for or-parallelism in this benchmark. However, they are *unprofitable* opportunities, because only one of the two branches will succeed, and the other will fail quickly. The trick here is to not get derailed by or-parallelism, so that the actual and-parallelism present (across the large tree) can be exploited.

The original Typed Racket implementation has a particular performance problem with `treecall`, making it much slower than the Haskell version⁸, so we omit a comparison between Racket and Haskell here. Rather, we evaluate several Haskell implementation variants. As described in §6, our idiomatic Haskell port is sequential, and our parallelism version is parameterized so as to take either

⁸The use of lists rather than sets in places was one factor that related in a severe slowdown. We’re still investigating the problems with the original Typed Racket version. But even the much faster Haskell version still takes long enough to get a parallel speedup!

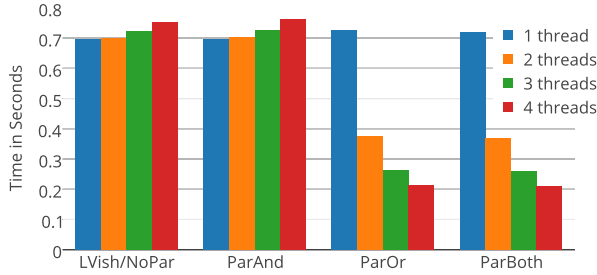


Figure 7. Bigcall(4): Highly polymorphic inputs benchmark on the desktop platform; shows time, in seconds, to type check (`map-with-funcs + - * /`), computing all solutions. Note that the time for Pure/Seq is too much larger than the LVish version to show here (7.93s for even the first solution).

a parallel or sequential implementation of conjunction and disjunction.

- Pure/Seq – original Racket to Haskell port.
- LVish/NoPar – refactored to use the LVish Par monad, and to use generators to represent solution streams. Still sequential.
- ParAnd – turn on parallel And.
- ParOr – turn on parallel Or.
- ParBoth – turn on parallel And and parallel Or.

As we mentioned, treecall provides unprofitable Or-parallelism with very poor granularity. For this reason the `ParOr` variant of the benchmark not only fails to speed up, but it gets a parallel slowdown due to useless tasks and threads bouncing between cores. Not only does the `ParAnd` variant work well, but `ParBoth` is fine as well. `ParAnd` can cure the problem with or-parallelism in this case, because the `andPar` function is essentially the “outer loop”, and it is this level of tasks that are stolen most in the underlying work-stealing implementation of the Par monad.

In summary, `Treecall(16)` takes 1.29s in the pure (non-LVish) Haskell version on the server platform. Then it gets up to a $7.68\times$ parallel speedup on the server platform when switching to LVish, and $3.17\times$ speedup on the desktop platform.

Bigcall results

Next, we evaluate `bigcall(3)` and `bigcall(4)`—three and four argument variants of `map-with-funcs`. The numbers we report here for LVish compute *all solutions*. For `bigcall` there are a modest number of solutions (hundreds) for very large search spaces (millions).

There are two distinct opportunities for performance in this benchmark:

- Deforestation: the Pure/Seq version uses lists to represent streams, and exploring the search space requires a lot of list book-keeping.
- Or-parallelism: because of the combinatorial explosion of different possible types for `(+)`, `(*)` and so on, there should be plenty of parallel work in exploring this search space in parallel.

Indeed, the LVish version of the program achieves a big speedup in both these categories. For example, the Pure/Seq Haskell implementation takes 2.36 seconds to compute all solutions for `bigcall(3)` (server platform). Just by deforesting intermediate lists using generators, LVish achieves a $9.37\times$ speedup over this baseline. And, even though the remaining time is only a short 0.25s, LVish `ParOr`

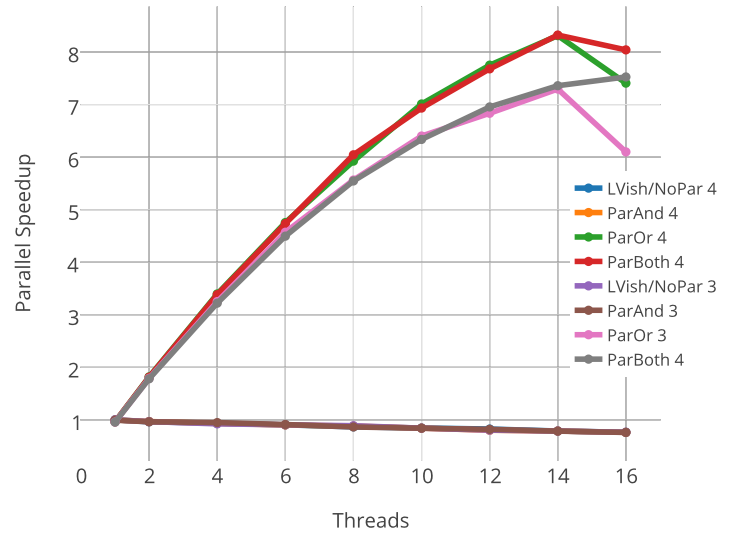


Figure 8. Bigcall(3) and bigcall(4): Parallel speedup on server platform, three- and four-argument version. Even bigcall(3) is time-consuming enough to get a good parallel speedup. This graph also demonstrates the “last core slowdown” which is still a problem with some parallel programs on GHC 7.8.3.

and `ParBoth` variants still achieve greater than $7\times$ parallel speedup, resulting in a total speedups of $70\times$ or higher over the original, idiomatic Haskell version ported from Racket. The speedup compared to original Racket version would be even greater, because the Racket version takes 3.0s to compute the first solution to `bigcall(3)`.

Note that in `bigcall`, computing all solutions is wasteful, with the purely functional Haskell version (Pure/Seq) taking only 0.18s as opposed to 2.36s on `bigcall(3)` to compute the first rather than all answers on the server platform. However, first, the deforestation benefit from switching to LVish makes up for the difference. And, second, even the Pure/Seq Haskell version is faster than the original Typed Racket implementation, with the original Typed Racket version taking 3.0s to compute just the *first* solution to (`map-with-funcs + - *`), i.e., `bigcall(3)`.

When scaling to the four-argument version, `bigcall(4)`, LVish is much faster than the purely functional versions even though it is computing all solutions. In fact, this example produces only a few hundred solutions, which is small compared to the size of the search space. Union types in Typed Racket programs represent or-parallelism with a very low *survival rate*—often only one of the variants in a union matches. For this larger case, the parallel speedups go to $8.46\times$ (server) and $3.43\times$ (desktop), and total speedup of LVish over Pure/Seq approaches $80\times$.

If we are willing to admit nondeterminism, it is extremely straightforward to have the parallel LVish implementation asynchronously report the first result it finds, and kill the rest of the `runPar` session. However, this is counter to our goal of deterministic parallelism. In future work, we plan to study the issue of extracting a *deterministic* result, in parallel.

7. Related Work

The *monad-par* system predated LVish and provided IVars as the sole synchronization construct. Marlow et al. [23] presents parallel type inference as a motivating example, but no implementa-

tion was evaluated. In fact, any implementation would have been severely limited due to the restriction that LVars only be written once. As a result, a given type variable could be constrained a *single* time, which is not compatible with most type systems (including Hindley-Milner).

Work on parallel Prolog [6, 13] solves similar issues of And-and Or-parallelism. Prolog’s control model is distinct from our LVish implementation, though, due to the presence of nondeterminism and the cut predicate. As a result, the data structures used are specific to logic programming, where ours are more general.

Saraswat and Rinard [27] discuss $cc(\downarrow, \rightarrow)$, a language for concurrent constraint logic programming. Their system also includes a notion of blocking reads analogous to the threshold reads of LVars, and additionally requires that concurrently written constraints be consistent with one another. It has nondeterministic operations, but also does identify a subset of operations that retain determinism. Modern concurrent constraint programming systems are available in software packages like Gecode [28]—a C++ library. The builtin constraint types and search strategies would not apply to, e.g., Typed Racket typechecking. But one could use such a system as an alternate starting point for this research: writig new C++ code to extend the system with new models, propagators, and branchers to handle the type checking. However, the verification or testing burden to ensure this C++ code retains determinism is a much more difficult obligation than with LVish. In LVish, we must ensure *only* that a TyVar LVar’s join function has the appropriate properties (commutative, associative, idempotence, and absorption). But join is just a pure function that can be tested for these properties with QuickCheck or other methods.

Deterministic Search Algorithms In [16], the authors give a deterministic parallel algorithm for backtracking search problems. COMMON-CRCW, their computational model, allows for arbitrary concurrent reads, and restricts concurrent writes by requiring that all threads write the same value.

IBM’s CPLEX system [7] offers a parallel solver for integer linear programs, with some extensions. Their solver is deterministic, except in the case where the user provides *control callbacks*, which allow observation and modification of the state of the parallel search.

8. Conclusion

Type checking in modern type systems is an expensive process, but not one that has previously been parallelized. We saw how the LVar framework is one possible way to address this challenge while also ensuring determinism in addition to gaining parallelism. We showed substantial parallel scaling and improvement in wall-clock time on two very different type systems: one very widely used, with $3.57\times$ parallel speedup, and the other slow and sharply in need of parallelization, with up to $4.71\times$ and $8.46\times$ parallel-speedup on our two benchmarks, respectively. We hope to extend our approach to accommodate choosing a single (first) answer deterministically, and apply our techniques in the original Typed Racket implementation.

References

- [1] L. Allison. Continuations implement generators and streams. *The Computer Journal*, 33(5):460–465, 1990. . URL <http://comjnl.oxfordjournals.org/content/33/5/460.abstract>.
- [2] Arvind, Rishiyur S. Nikhil, and Keshav K. Pingali. I-structures: Data structures for parallel computing. *ACM Trans. Program. Lang. Syst.*, 11:598–632, October 1989. ISSN 0164-0925. . URL <http://doi.acm.org/10.1145/69558.69562>.
- [3] Robert L Bocchino Jr, Vikram S Adve, Danny Dig, Sarita V Adve, Stephen Heumann, Rakesh Komuravelli, Jeffrey Overbey, Patrick Simmons, Hyojin Sung, and Mohsen Vakilian. A type and effect system for deterministic parallel java. *ACM Sigplan Notices*, 44(10):97–116, 2009.
- [4] Edwin Brady. Idris, a general-purpose dependently typed programming language: Design and implementation. *Journal of Functional Programming*, 23(05):552–593, 2013.
- [5] Koen Claessen, Mary Sheeran, and Bo Joel Svensson. Expressive array constructs in an embedded gpu kernel programming language. In *Proceedings of the 7th workshop on Declarative aspects and applications of multicore programming*, pages 21–30. ACM, 2012.
- [6] Vitor Santos Costa, David HD Warren, and Rong Yang. *Andorra I: a parallel Prolog system that transparently exploits both And-and or-parallelism*, volume 26. ACM, 1991.
- [7] IBM ILOG CPLEX. V12. 1: Users manual for cplex. *International Business Machines Corporation*, 46(53):157, 2009.
- [8] Luis Damas and Robin Milner. Principal type-schemes for functional programs. In *Proceedings of the 9th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 207–212. ACM, 1982.
- [9] Eelco Dolstra and Andres Löh. Nixos: A purely functional linux distribution. In *ACM Sigplan Notices*, volume 43, pages 367–378. ACM, 2008.
- [10] Matthew Flatt and PLT. Reference: Racket. Technical report, PLT Design, Inc., 2010. <http://racket-lang.org/tr1/>.
- [11] Ian P Gent, Chris Jefferson, Ian Miguel, NC Moore, Peter Nightingale, Patrick Prosser, and Chris Unsworth. A preliminary review of literature on parallel constraint solving. In *Proceedings PMCS 2011 Workshop on Parallel Methods for Constraint Solving*. Citeseer, 2011.
- [12] Ralph E. Griswold and Madge T. Griswold. History of programming languages—ii. chapter History of the Icon Programming Language, pages 599–624. ACM, New York, NY, USA, 1996. ISBN 0-201-89502-1. . URL <http://doi.acm.org/10.1145/234286.1057830>.
- [13] Gopal Gupta and Vitor Santos Costa. And-or parallelism in full prolog with paged binding arrays. In *PARLE’92 Parallel Architectures and Languages Europe*, pages 617–632. Springer, 1992.
- [14] Youssef Hamadi, Said Jabbour, and Lakhdar Sais. Manysat: a parallel sat solver. *JSAT*, 6(4):245–262, 2009.
- [15] Youssef Hamadi, Said Jabbour, and Jabbour Sais. Control-based clause sharing in parallel sat solving. In *Autonomous Search*, pages 245–267. Springer, 2012.
- [16] Kieran T Herley, Andrea Pietracaprina, and Geppino Pucci. Deterministic parallel backtrack search. *Theoretical Computer Science*, 270(1): 309–324, 2002.
- [17] Ranjit Jhala. Refinement types for haskell. In *Proceedings of the ACM SIGPLAN 2014 workshop on Programming languages meets program verification*, pages 27–28. ACM, 2014.
- [18] Lindsey Kuper and Ryan R Newton. Lvars: lattice-based data structures for deterministic parallelism. In *Proceedings of the 2nd ACM SIGPLAN workshop on Functional high-performance computing (FHPC’13)*, pages 71–84. ACM, 2013.
- [19] Lindsey Kuper, Aaron Todd, Sam Tobin-Hochstadt, and Ryan R Newton. Taming the parallel effect zoo: Extensible deterministic parallelism with lvish. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, page 2. ACM, 2014.
- [20] Lindsey Kuper, Aaron Turon, Neelakantan R. Krishnaswami, and Ryan R. Newton. Freeze after writing: quasi-deterministic parallel programming with lvars. In *POPL*, pages 257–270, 2014.
- [21] John Launchbury and Simon L Peyton Jones. Lazy functional state threads. In *ACM SIGPLAN Notices*, volume 29, pages 24–35. ACM, 1994.
- [22] Barbara Liskov. A history of clu. *SIGPLAN Not.*, 28(3):133–147, March 1993. ISSN 0362-1340. . URL <http://doi.acm.org/10.1145/155360.155367>.

- [23] Simon Marlow, Ryan Newton, and Simon Peyton Jones. A monad for deterministic parallelism. In *Proceedings of the 4th ACM symposium on Haskell*, Haskell '11, pages 71–82. ACM, 2011.
- [24] Simon Peyton Jones, Alastair Reid, Fergus Henderson, Tony Hoare, and Simon Marlow. A semantics for imprecise exceptions. In *ACM SIGPLAN Notices*, volume 34, pages 25–36. ACM, 1999.
- [25] Benjamin C. Pierce and David N. Turner. Local Type Inference. *ACM Trans. Program. Lang. Syst.*, 22(1):1–44, 2000.
- [26] Tarun Prabhu, Shreyas Ramalingam, Matthew Might, and Mary Hall. Eigencfa: accelerating flow analysis with gpus. In *ACM SIGPLAN Notices*, volume 46, pages 511–522. ACM, 2011.
- [27] Vijay A Saraswat and Martin Rinard. Concurrent constraint programming. In *Proceedings of the 17th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 232–245. ACM, 1989.
- [28] Christian Schulte and Guido Tack. Views and iterators for generic constraint implementations. In *Recent Advances in Constraints (2005)*, volume 3978 of *Lecture Notes in Artificial Intelligence*, pages 118–132. Springer-Verlag, 2006. URL <http://www.gecode.org/paper.html?id=SchulteTack:Advances:2006>.
- [29] Mary Shaw, William A. Wulf, and Ralph L. London. Abstraction and verification in alphard: Defining and specifying iteration and generators. *Commun. ACM*, 20(8):553–564, August 1977. ISSN 0001-0782. . URL <http://doi.acm.org/10.1145/359763.359782>.
- [30] Jeremy G. Siek and Walid Taha. Gradual typing for functional languages. In *Scheme and Functional Programming Workshop, University of Chicago, Technical Report TR-2006-06*, pages 81–92, September 2006.
- [31] Jeremy G Siek and Walid Taha. Gradual typing for functional languages. In *Scheme and Functional Programming Workshop*, volume 6, pages 81–92, 2006.
- [32] Vincent St-Amour, Sam Tobin-Hochstadt, Matthew Flatt, and Matthias Felleisen. Typing the numeric tower. pages 289–303, 2012.
- [33] T. Stephen Strickland, Sam Tobin-Hochstadt, , and Matthias Felleisen. Practical variable-arity polymorphism. pages 32–46, March 2009.
- [34] Sam Tobin-Hochstadt and Matthias Felleisen. Interlanguage migration: from scripts to programs. pages 964–974. (Companion (Dynamic Languages Symposium), 2006.
- [35] Sam Tobin-Hochstadt and Matthias Felleisen. The design and implementation of Typed Scheme. pages 395–406, 2008.
- [36] Sam Tobin-Hochstadt and Matthias Felleisen. Logical types for untyped languages. In *Proceedings of the 15th ACM SIGPLAN International Conference on Functional Programming, ICFP '10*, pages 117–128, New York, NY, USA, 2010. ACM. ISBN 978-1-60558-794-3. . URL <http://doi.acm.org/10.1145/1863543.1863561>.
- [37] Pascal Van Hentenryck, Vijay Saraswat, and Yves Deville. Design, implementation, and evaluation of the constraint language cc (fd). *The Journal of Logic Programming*, 37(1-3):139–164, 1998.
- [38] A. Zobel. Program structure as basis for parallelizing global register allocation. In *Computer Languages, 1992., Proceedings of the 1992 International Conference on*, pages 262–271, Apr 1992. .