

Flow-Sensitive Type Qualifiers*

Jeffrey S. Foster

Tachio Terauchi

Alex Aiken

EECS Department
University of California, Berkeley
Berkeley, CA 94720-1776
{jfoster,tachio,aiken}@cs.berkeley.edu

ABSTRACT

We present a system for extending standard type systems with flow-sensitive type qualifiers. Users annotate their programs with type qualifiers, and inference checks that the annotations are correct. In our system only the type qualifiers are modeled flow-sensitively—the underlying standard types are unchanged, which allows us to obtain an efficient constraint-based inference algorithm that integrates flow-insensitive alias analysis, effect inference, and ideas from linear type systems to support strong updates. We demonstrate the usefulness of flow-sensitive type qualifiers by finding a number of new locking bugs in the Linux kernel.

Categories and Subject Descriptors

D.2.1 [Software Engineering]: Requirements/Specifications; D.2.4 [Software Engineering]: Software/Program Verification; D.3.3 [Programming Languages]: Language Constructs and Features; F.3.1 [Logics and Meanings of Programs]: Specifying and Verifying and Reasoning about Programs; F.3.3 [Logics and Meanings of Programs]: Studies of Program Constructs

General Terms

Algorithms, Design, Reliability, Experimentation, Languages, Theory, Verification

Keywords

Types, type qualifiers, alias analysis, effect inference, flow-sensitivity, constraints, restrict, locking, Linux kernel

*This research was supported in part by NSF CCR-9457812, NASA Contract No. NAG2-1210, NSF CCR-0085949, and DARPA Contract No. F33615-00-C-1693.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PLDI'02, June 17-19, 2002, Berlin, Germany.

Copyright 2002 ACM 1-58113-463-0/02/0006 ...\$5.00.

1. INTRODUCTION

Standard type systems are *flow-insensitive*, meaning a value's type is the same everywhere. However, many important program properties are *flow-sensitive*. Checking such properties requires associating different facts with a value at different program points.

This paper shows how to extend standard type systems with user-specified flow-sensitive *type qualifiers*, which are atomic properties that refine standard types. In our system users annotate programs with type qualifiers, and inference checks that the annotations are correct. The critical feature of our approach is that flow-sensitivity is restricted to the type qualifiers that decorate types—the underlying standard types are unchanged—which allows us to obtain an efficient type inference algorithm. Type qualifiers capture a natural class of flow-sensitive properties, while efficient inference of the type qualifiers allows us to apply an implementation to large code bases with few user annotations.

As an example of type qualifiers, consider the type `File` used for I/O operations on files. In most systems `File` operations can only be used in certain ways: a file must be opened for reading before it is read, it must be opened for writing before it is written to, and once closed a file cannot be accessed. We can express these rules with flow-sensitive type qualifiers. We introduce qualifiers `open`, `read`, `write`, `readwrite`, and `closed`. The type `open File` describes a file that has been opened in an unknown mode, the type `read File` (respectively `write File`) is a file that is open for reading (respectively writing), the type `readwrite File` is a file open for both reading and writing, and the type `closed File` is a closed file. These qualifiers capture inherently flow-sensitive properties. For example, the `close()` function takes an `open File` as an argument and changes the file's state to `closed File`.

These five qualifiers have a natural subtyping relation: `readwrite` \leq `read` \leq `open` and `readwrite` \leq `write` \leq `open`. The qualifier `closed` is incomparable to other qualifiers because a file may not be both closed and open. Qualifiers that introduce subtyping are very common, and our framework supports subtyping directly; in addition to a set of qualifiers, users can define a partial order on the qualifiers.

Our results build on recent advances in flow-sensitive type systems [5, 7, 25] as well as our own previous work on flow-insensitive type qualifiers [16, 24]. The main contribution of our work is a practical, flow-sensitive type inference algorithm, in contrast to the type checking systems of [5, 7, 25].

Our flow-sensitive type inference algorithm is made prac-

tical by solving constraints lazily. As in any flow-sensitive analysis, explicitly forming a model of the store at every program point is prohibitively expensive for large code bases. By generating a constraint system linear in the size of the type-annotated program and solving only the portion of the constraints needed to check qualifier annotations, our algorithm is able to scale to large examples.

Finally, our system is designed to be sound; we aim to prove the absence of bugs, not just to be heuristically good at finding bugs. For example, we believe that our system could be integrated into Java in a sound manner. We have shown soundness for `restrict` (Section 4), a key new construct in our system (see technical report [15]). Since the remainder of our system can be viewed as a simplification of [25], we believe it is straightforward to prove soundness for our full type system using their techniques.

In Section 5 we report on experience with two applications, analyzing locking behavior in the Linux kernel and analyzing C stream library usage in application code. Our system found a number of new locking bugs, including some that extend across multiple functions or even, in one case, across multiple files.

1.1 System Architecture

Our flow-sensitive qualifier inference algorithm has several interlocking components. We first give an overview of the major pieces and how they fit together.

We expect programmers to interact with our type system, both when adding qualifier annotations and when reviewing the results of inference. Thus, we seek a system that supports efficient inference and is straightforward for a programmer to understand and use. Our type inference system integrates alias analysis, effect inference, and ideas from linear type systems.

- We use a flow-insensitive alias analysis to construct a model of the store. The alias analysis infers an *abstract location* for the result of each program expression; expressions that evaluate to the same abstract location may be aliased.
- We use effect inference [20] to calculate the set of abstract locations an expression e might use during e 's evaluation. These effects are used in analyzing function calls and `restrict` (see below). Effect inference is done simultaneously with alias analysis.
- We model the state at a program point as an abstract *store*, which is a mapping from abstract locations to types. We can use the abstract locations from the flow-insensitive alias analysis because we allow only the type qualifiers, and not the underlying standard types, to change during execution. We represent abstract stores using a constraint formalism. Store constructors model allocations, updates, and function calls, and store constraints $C_1 \leq C_2$ model a branch from the program point represented by store C_1 to the program point represented by store C_2 .
- We compute a linearity [25] for each abstract location at each program point. Informally, an abstract location is *linear* if the type system can prove that it corresponds to a single concrete location in every execution; otherwise, it is non-linear. We perform *strong updates*

[4] on locations that are linear and *weak updates* on locations that are non-linear. A strong update can change the qualifier on a location's type arbitrarily. Weak updates cannot change qualifiers. Computing linearities is important because most interesting flow-sensitive properties require strong updates.

- The system described so far has a serious practical weakness: Type inference may fail because a location on which a strong update is needed may be inferred to be non-linear. We address this with a new annotation `restrict`. The expression `restrict $x=e$ in e'` introduces a new name x bound to the value of e . The name x is given a fresh abstract location, and among all aliases of e , only x and values derived from x may be used within e' . Thus the location of x may be linear, and hence may be strongly updated, even if the location of e is non-linear. We use effects to enforce the correctness of `restrict` expressions—soundness requires that the location of e does not appear in the effect of e' .
- We use effects to increase the precision of the analysis. If an expression e does not reference location ρ , which we can determine by examining the effect of e , then it does not access the value stored at ρ , and the analysis of ρ can simply flow from the store preceding e to the one immediately after e without passing through e . If e is an application of a function called in many different contexts, then this idea makes e fully polymorphic in all the locations that e does not reference.

2. RELATED WORK

We discuss three threads of related work: type systems, dataflow analysis, and tools for finding bugs in software.

Type Systems. Our type system is inspired by region and alias type checking systems designed for low-level programs [5, 25, 29]. Two recent language proposals, Vault [7] and Cyclone [17], adapt similar ideas for checking high-level programs. Both of these languages are based on type checking and require programmers to annotate their programs with types. In contrast, we propose a simpler and less expressive monomorphic type system that is designed for efficient type inference. Our system incorporates effect inference [20, 32] to gain a measure of polymorphism. Recent work on Vault [12] includes a construct `focus` that is similar to `restrict`.

The type state system of NIL [27] is one of the earliest to incorporate flow-sensitive type checking. Xu et al [33] use a flow-sensitive analysis to check type safety of machine code. Type systems developed for Java byte code [22, 26] also incorporate flow-sensitivity to check for initialization before use and to allow reuse of the same local variable with different types.

Igarashi and Kobayashi [18] propose a general framework for resource usage analysis, which associates a trace with each object specifying valid accesses to the object and checks that the program satisfies the trace specifications. They provide an inference algorithm, although it is unclear how efficient it is in practice since it invokes as a sub-step an unspecified algorithm to check that a trace set is valid.

Flanagan and Freund [13] use a type checking system to verify Java locking behavior. In Java locks are acquired and

released according to a lexical discipline. To model locking in the Linux kernel (as in Section 5) we must allow non-lexically scoped lock acquires and releases.

The subset of our system consisting of alias analysis and effect inference can be seen as a monomorphic variant of region inference [28]. The improvements to region inference reported in [2] are a much more expensive and precise method for computing linearities.

Dataflow Analysis. Although our type-based approach is related to dataflow analysis [1], it differs from classical dataflow analysis in several ways. First, we generate constraints over stores and types to model the program. Thus there is no distinction between forward and backward analysis; information may flow in both directions during constraint resolution, depending on the specified qualifier partial order. Second, we explicitly handle pointers, heap-allocated data, aliasing, and strong/weak updates. Third, there is no distinction between interprocedural and intraprocedural analysis in our system.

The strong/weak update distinction was first described by Chase et al [4]. Several techniques that allow strong updates have been proposed for dataflow-based analysis of programs with pointers, among them [3, 8, 31]. Jagannathan et al [19] present a system for must-alias analysis of higher-order languages. The linearity computation in our system corresponds to their singleness computation, and they use a similar technique to gain polymorphism by flowing some bindings around function calls.

Another recent system for checking typestate properties is ESP [6]. Like our proposal, ESP incorporates a conservative alias analysis. There are also significant differences: ESP is more directly based on dataflow analysis and incorporates a path-sensitive symbolic execution component. ESP has been used to check the correctness of C stream library usage in gcc.

Bug-Finding Tools. The AST Toolkit provides a framework for posing user-specified queries on abstract syntax trees annotated with type information. The AST Toolkit has been successfully used to uncover many bugs [30].

Meta-level compilation [9] is a system for finding bugs in programs. The programmer specifies a flow-sensitive property as an finite state automaton. A program is analyzed by traversing control paths and triggering state transitions of the automata on particular actions in program statements. The system warns of potential errors when an automaton enters an error state. In [9] an intraprocedural analysis of lock usage in the Linux kernel uncovered many local locking bugs. Our type-based system found *interprocedural* locking bugs that extended across multiple functions or even, in one case, across multiple files (Section 5).¹ Newer work on meta-level compilation [10] includes some interprocedural dataflow, but it is unclear how their interprocedural dataflow analysis handles aliasing.

LCLint [11] is a dataflow-based tool for checking properties of programs. To use LCLint, the programmer adds extra annotations to their program. LCLint performs flow-sensitive intraprocedural analysis, using the programmer’s

¹The bugs were found in a newer version of the Linux kernel than examined by [9], so a direct comparison is not possible, though these bugs cannot be found by purely intraprocedural analysis.

annotations at function calls.

ESC/Java [14] is a tool for finding errors in Java programs. ESC/Java uses sophisticated theorem-proving technology to verify program properties, and it includes a rich language for program annotations.

3. TYPE SYSTEM

We describe our type system using a call-by-value lambda calculus extended with pointers and type qualifier annotations. The source language is

$$e ::= x \mid n \mid \lambda x.e \mid e_1 e_2 \mid \mathbf{ref} e \mid !e \mid e_1 := e_2 \mid \mathbf{assert}(e, Q) \mid \mathbf{check}(e, Q)$$

Here x is a variable, n is an integer, $\lambda x.e$ is a function with argument x and body e , the expression $e_1 e_2$ is the application of function e_1 to argument e_2 , the expression $\mathbf{ref} e$ allocates memory and initializes it to e , the expression $!e$ dereferences pointer e , and the expression $e_1 := e_2$ assigns the value of e_2 to the location e_1 points to.

We introduce qualifiers into the source language by adding two new forms [16]. The expression $\mathbf{assert}(e, Q)$ asserts that e ’s top-level qualifier is Q , and the expression $\mathbf{check}(e, Q)$ type checks only if e ’s top-level qualifier is at most Q .

Our type inference algorithm is divided into two steps. First we perform an initial flow-insensitive alias analysis and effect inference. Second we generate and solve store and qualifier constraints and compute linearities.

3.1 Alias Analysis and Effect Inference

We present the flow-insensitive alias analysis and effect inference as a translation system rewriting source expressions to expressions decorated with locations, types, and effects. The target language is

$$\begin{aligned} e &::= x \mid n \mid \lambda^L x:t.e \mid e_1 e_2 \mid \mathbf{ref}^\rho e \mid !e \mid e_1 := e_2 \\ &\quad \mid \mathbf{assert}(e, Q) \mid \mathbf{check}(e, Q) \\ t &::= \alpha \mid \mathit{int} \mid \mathit{ref}(\rho) \mid t \xrightarrow{L} t' \\ L &::= \psi \mid \{\rho\} \mid L_1 \cup L_2 \mid L_1 \cap L_2 \end{aligned}$$

The target language extends the source language syntax in two ways. Every allocation site $\mathbf{ref}^\rho e$ is annotated with the abstract location ρ that is allocated, and each function $\lambda^L x:t.e$ is annotated with both the type t of its parameter and the effect L of calling the function. Effects are unions and intersections of *effect variables* ψ , which represent an unknown set of effects that effect inference solves for, and *effect constants* ρ , which stands for either a read, write, or allocation of location ρ . For simplicity in this paper we do not distinguish which of the three possible effects ρ stands for, although we do so in our implementation.

Foreshadowing flow-sensitive analysis, pointer types are written $\mathit{ref}(\rho)$, and we maintain a separate global abstract store C_I mapping locations ρ to types; $C_I(\rho) = \tau$ if location ρ contains data of type τ . If type inference requires $\rho = \rho'$, we also require $C_I(\rho) = C_I(\rho')$. Function types $t \xrightarrow{L} t'$ contain the effect L of calling the function.

Figure 1 gives rules for performing alias analysis and effect inference while translating source programs into our target language. This translation system proves judgments $\Gamma \vdash e \Rightarrow e' : t; L$, meaning that in type environment Γ , expression e translates to expression e' , which has type t , and the evaluation of e may have effect L .

$$\begin{array}{c}
\frac{x \in \text{dom}(\Gamma)}{\Gamma \vdash x \Rightarrow x : \Gamma(x); \emptyset} \text{ (Var)} \\
\\
\frac{}{\Gamma \vdash n \Rightarrow n : \text{int}; \emptyset} \text{ (Int)} \\
\\
\frac{\Gamma \vdash e \Rightarrow e' : t; L \quad C_I(\rho) = t \quad \rho \text{ fresh}}{\Gamma \vdash \text{ref } e \Rightarrow \text{ref}^\rho e' : \text{ref}(\rho); L \cup \{\rho\}} \text{ (Ref)} \\
\\
\frac{\Gamma \vdash e \Rightarrow e' : t; L \quad t = \text{ref}(\rho) \quad \rho \text{ fresh}}{\Gamma \vdash !e \Rightarrow !e' : C_I(\rho); L \cup \{\rho\}} \text{ (Deref)} \\
\\
\frac{\Gamma \vdash e_1 \Rightarrow e'_1 : t_1; L_1 \quad \Gamma \vdash e_2 \Rightarrow e'_2 : t_2; L_2 \quad t_1 = \text{ref}(\rho) \quad C_I(\rho) = t_2 \quad \rho \text{ fresh}}{\Gamma \vdash e_1 := e_2 \Rightarrow e'_1 := e'_2 : t_2; L_1 \cup L_2 \cup \{\rho\}} \text{ (Assign)} \\
\\
\frac{\Gamma[x \mapsto \alpha] \vdash e \Rightarrow e' : t; L \quad L \subseteq \psi \quad \alpha, \psi \text{ fresh}}{\Gamma \vdash \lambda x. e \Rightarrow \lambda^\psi x. \alpha. e' : \alpha \rightarrow^\psi t; \emptyset} \text{ (Lam)} \\
\\
\frac{\Gamma \vdash e_1 \Rightarrow e'_1 : t_1; L_1 \quad \Gamma \vdash e_2 \Rightarrow e'_2 : t_2; L_2 \quad t_1 = t_2 \xrightarrow{\psi} \beta \quad \psi, \beta \text{ fresh}}{\Gamma \vdash e_1 e_2 \Rightarrow e'_1 e'_2 : \beta; L_1 \cup L_2 \cup \psi} \text{ (App)} \\
\\
\frac{\Gamma \vdash e \Rightarrow e' : t; L}{\Gamma \vdash \text{assert}(e, Q) \Rightarrow \text{assert}(e', Q) : t; L} \text{ (Assert)} \\
\\
\frac{\Gamma \vdash e \Rightarrow e' : t; L}{\Gamma \vdash \text{check}(e, Q) \Rightarrow \text{check}(e', Q) : t; L} \text{ (Check)} \\
\\
\frac{\Gamma \vdash e \Rightarrow e' : t; L}{\Gamma \vdash e \Rightarrow e' : t; L \cap (\text{locs}(\Gamma) \cup \text{locs}(t))} \text{ (Down)}
\end{array}$$

Figure 1: Type, alias, and effect inference

The set of locations appearing in a type, $\text{locs}(t)$, is

$$\begin{aligned}
\text{locs}(\text{int}) &= \emptyset \\
\text{locs}(\text{ref}(\rho)) &= \{\rho\} \cup \text{locs}(C_I(\rho)) \\
\text{locs}(t_1 \xrightarrow{L} t_2) &= \text{locs}(t_1) \cup \text{locs}(t_2) \cup L
\end{aligned}$$

We assume that $\text{locs}(\alpha)$ is empty until α is equated with a constructed type. We define $\text{locs}(\Gamma)$ to be $\bigcup_{[x \mapsto t \in \Gamma]} \text{locs}(t)$.

We briefly discuss the rules in Figure 1:

- (Var) and (Int) are standard. In lambda calculus, a variable is an r -value, not an l -value, and accessing a variable has no effect.
- (Ref) allocates a fresh abstract location ρ . We add the effect $\{\rho\}$ of the allocation to the effect and record in C_I the type to which the location ρ points.
- (Deref) evaluates e , which yields a value of type t . As is standard in type inference, to compute the location e points to we create a fresh location ρ and equate the type t with the type $\text{ref}(\rho)$. We look up the type of location ρ in C_I and add ρ to the effect set.
- (Assign) writes a location. Note that the type of e_2 and the type that e_1 points to are equated. Because types contain locations, this forces potentially aliased locations to be modeled by one abstract location.
- (Lam) defines a function. We annotate the function with the effect ψ of the function body and the type α of the parameter. Function types always have an effect

<pre> fun f w = let x = ref 0 y = ref(assert(1, q_a)) z = ref(assert(2, q_b)) in /* Write to x's cell */ x := 3; w := 4; y := assert(5, q_c); if (...) f z; check(!y, q_c) </pre>	<pre> fun^{ρ_z} f w: ref(ρ_z) = let x = ref^{ρ_x} 0 y = ref^{ρ_y}(assert(1, q_a)) z = ref^{ρ_z}(assert(2, q_b)) in x := 3; w := 4; y := assert(5, q_c); if (...) f z; check(!y, q_c) </pre>
--	---

(a) Source program

(b) Target program

$$C_I(\rho_x) = C_I(\rho_y) = C_I(\rho_z) = \text{int}$$

Figure 2: Example alias and effect analysis

variable ψ on the arrow, which makes effect inference easier. Notice that creating a function has no effect (the potential allocation of a closure does not count as an effect, because a closure cannot be updated).

- (App) applies a function to an argument. The effect of applying e_1 to e_2 includes the effect ψ of calling the function e_1 represents. Notice that e_1 's argument type is constrained to be equal to the type of e_2 . As before, this forces possibly-aliased locations to have the same abstract location.
- (Assert) and (Check) are translated unchanged into the target language. Qualifiers are flow-sensitive, so we do not model them during this first, flow-insensitive step of the algorithm.
- (Down) hides effects on purely local state. If evaluating e produces an effect on some location ρ neither in Γ nor in t , then ρ cannot be accessed in subsequent computation. Thus we can conservatively approximate the set of effects that may be visible as $\text{locs}(\Gamma) \cup \text{locs}(t)$. By intersecting the effects L with the set of effects that may be visible, we increase the precision of effect inference, which in turn increases the precision of flow-sensitive type qualifier inference. Although (Down) is not a syntactic rule, it only needs to be applied once per function body [15].

Figure 2 shows an example program and its translation. We use some syntactic sugar; all of these constructs can be encoded in our language (e.g., by assuming a primitive **Y** combinator of the appropriate type). In this example the constant qualifiers q_a , q_b , and q_c are in the discrete partial order (the qualifiers are incomparable). Just before f returns, we wish to check that y has the qualifier q_c . This check succeeds only if we can model the update to y as a strong update.

In Figure 2, we assign x , y , and z distinct locations ρ_x , ρ_y , and ρ_z , respectively. Because f is called with argument z and our system is not polymorphic in locations, our alias analysis requires that the types of z and w match, and thus w is given the type $\text{ref}(\rho_z)$. Finally, notice that since x and y are purely local to the body of f , using the rule (Down) our analysis hides all effects on ρ_x and ρ_y . The effect of f is $\{\rho_z\}$ because f writes to its parameter w , which has type $\text{ref}(\rho_z)$. (More precisely, f has effect ψ where $\{\rho_z\} \subseteq \psi$.)

Let n be the size of the input program. Applying the rules in Figure 1 generates a constraint system of size $O(n)$, using a suitable representation of $\text{locs}(\Gamma) \cup \text{locs}(t)$ (see [15]). Resolving the type equality constraints in the usual way with unification takes $O(n\alpha(n))$ time, where $\alpha(\cdot)$ is the inverse Ackerman’s function. The remaining constraints are *effect constraints* of the form $L \subseteq \psi$. We solve these constraints on-demand—in the next step of the algorithm we ask queries of the form $\rho \in L$. We can answer all such queries for a single location ρ in $O(n)$ time [15].

3.2 Stores and Qualified Types

Next we perform flow-sensitive analysis to check the qualifier-related annotations. In this second step of the algorithm we take as input a program that has been decorated with types, locations, and effects by the inference algorithm of Figure 1. Throughout this step we treat the abstract locations ρ and effects L from the first step as constants. We analyze the input program using the extended types shown below:

$$\begin{aligned} \tau &::= Q \sigma \\ Q &::= \kappa \mid B \\ \sigma &::= \alpha \mid \text{int} \mid \text{ref}(\rho) \mid (C, \tau) \longrightarrow^L (C', \tau') \\ C &::= \varepsilon \mid \text{Alloc}(C, \rho) \mid \text{Assign}(C, \rho : \tau) \\ &\quad \mid \text{Merge}(C, C', L) \mid \text{Filter}(C, L) \\ \eta &::= 0 \mid 1 \mid \omega \end{aligned}$$

Here *qualified types* τ are standard types with qualifiers inserted at every level. Qualifiers Q are either *qualifier variables* κ , which stand for currently unknown qualifiers, or *constant qualifiers* B , specified by the user. We assume a supplied partial order \leq among constant qualifiers.

The flow-sensitive analysis associates a *store* C with each program point. This is in contrast to the flow-insensitive step, which uses one global store C_I to give types to locations. Function types are extended to $(C, \tau) \longrightarrow^L (C', \tau')$, where C describes the store the function is invoked in and C' describes the store when the function returns.

Each location in each store has an associated *linearity* η . There are three linearities: 0 for unallocated locations, 1 for linear locations (these admit strong updates), and ω for non-linear locations (which admit only weak updates). The three linearities form a lattice $0 < 1 < \omega$. Addition on linearities is as expected: $0 + x = x$, $1 + 1 = \omega$, and $\omega + x = \omega$.

A store is a vector

$$\{\rho_1^{\eta_1} : \tau_1, \dots, \rho_n^{\eta_n} : \tau_n\}$$

that assigns a type τ_i and a linearity η_i to every abstract location ρ_i computed by the alias analysis. We call such a vector a *ground store*. If G is a ground store, we write $G(\rho)$ for ρ ’s type in G , and we write $G_{\text{lin}}(\rho)$ for ρ ’s linearity in G .

Rather than explicitly associating a ground store with every program point, we represent stores using a constraint formalism. As the base case, we model an unknown store using a *store variable* ε . We relate stores at consecutive program points either with *store constructors* (see below), which build new stores from old stores, or with *store constraints* $C_1 \leq C_2$, which are generated at branches from the program point represented by store C_1 to the program point represented by store C_2 .

A *solution* to a system of store constraints is a mapping from store variables to ground stores, and from $\text{Assign}(\dots)$ stores (see below) to types. A solution S *satisfies* a system of store constraints if for each constraint $C_1 \leq C_2$ we have

$$\begin{aligned} &\frac{Q \leq Q'}{Q \text{ int} \leq Q' \text{ int}} \text{ (Int}_{\leq}) \\ &\frac{Q \leq Q'}{Q \text{ ref}(\rho) \leq Q' \text{ ref}(\rho)} \text{ (Ref}_{\leq}) \\ &\frac{Q \leq Q' \quad \tau_2 \leq \tau_1 \quad \tau'_1 \leq \tau'_2 \quad C_2 \leq C_1 \quad C'_1 \leq C'_2}{Q(C_1, \tau_1) \longrightarrow^L (C'_1, \tau'_1) \leq Q'(C_2, \tau_2) \longrightarrow^L (C'_2, \tau'_2)} \text{ (Fun}_{\leq}) \\ &\frac{\tau_i \leq \tau'_i \quad \eta_i \leq \eta'_i \quad i = 1..n}{\{\rho_1^{\eta_1} : \tau_1, \dots, \rho_n^{\eta_n} : \tau_n\} \leq \{\rho_1^{\eta'_1} : \tau'_1, \dots, \rho_n^{\eta'_n} : \tau'_n\}} \text{ (Store}_{\leq}) \end{aligned}$$

Figure 3: Store compatibility rules

$S(C_1) \leq S(C_2)$ according to the rules in Figure 3, and the solution satisfies the rules in Figure 4.

In Figure 3, constraints between stores yield constraints between linearities and types, which in turn yield constraints between qualifiers and between stores. In our constraint resolution algorithm, we exploit the fact that we are only interested in qualifier relationships to solve as little of the expensive store constraints as possible (see Section 3.4).

In (Ref_{\leq}) we require that the locations on the left- and right-hand sides of the \leq are the same. Alias analysis enforces this property, which corresponds to the standard requirement that subtyping becomes equality below a pointer constructor. We emphasize that in this step we treat abstract locations ρ as constants, and we will never attempt (or need) to unify two distinct locations to satisfy (Ref_{\leq}) .

In (Fun_{\leq}) we require that the effects of the constrained function types match exactly. It would also be sound to allow the effect of the left-hand function to be a subset of the effect of the right-hand function.

Figure 4 formalizes the four kinds of store constructors by showing how a solution S behaves on constructed stores.

The store $\text{Alloc}(C, \rho)$ is the same as store C , except that location ρ has been allocated once more. Allocating location ρ does not affect the types in the store but increases the linearity of location ρ by one.

The store $\text{Merge}(C, C', L)$ combines stores C and C' according to effect L . If $\rho \in L$, then $\text{Merge}(C, C', L)$ assigns ρ the type it has in C , otherwise $\text{Merge}(C, C', L)$ assigns ρ the type it has in C' . The linearity definition is similar.

The store $\text{Filter}(C, L)$ assigns the same types and linearities as C for all locations ρ such that $\rho \in L$. The types of all other locations are undefined, and the linearities of all other locations are 0.

Finally, the store $\text{Assign}(C, \rho : \tau)$ is the same as store C , except location ρ has been updated to type τ' where $\tau \leq \tau'$ (we allow a subtyping step here). If ρ is non-linear in C , then in Figure 4(c) we require that the type of ρ in $\text{Assign}(C, \rho : \tau)$ be at least the type of ρ in C ; this corresponds to a weak update. (In our implementation we require equality here.) Putting these together, intuitively if ρ is linear then its type in $\text{Assign}(C, \rho : \tau)$ is τ , otherwise its type is $\tau \sqcup S(C)(\rho)$, where \sqcup is the least-upper bound.

3.3 Flow-Sensitive Constraint Generation

Figure 5 gives the type inference rules for our system. In this system judgments have the form $\Gamma, C \vdash e : \tau, C'$, meaning that in type environment Γ and with initial store

$$\begin{aligned}
S(\text{Alloc}(C, \rho'))(\rho) &= S(C)(\rho) \\
S(\text{Merge}(C, C', L))(\rho) &= \begin{cases} S(C)(\rho) & \rho \in L \\ S(C')(\rho) & \text{otherwise} \end{cases} \\
S(\text{Filter}(C, L))(\rho) &= S(C)(\rho) \quad \rho \in L \\
S(\text{Assign}(C, \rho' : \tau))(\rho) &= \begin{cases} \tau' & \text{where } \tau \leq \tau' \quad \rho = \rho' \\ S(C)(\rho) & \text{otherwise} \end{cases}
\end{aligned}$$

(a) Types

$$\begin{aligned}
S(\text{Alloc}(C, \rho'))_{\text{lin}}(\rho) &= \begin{cases} 1 + S(C)_{\text{lin}}(\rho) & \rho = \rho' \\ S(C)_{\text{lin}}(\rho) & \text{otherwise} \end{cases} \\
S(\text{Merge}(C, C', L))_{\text{lin}}(\rho) &= \begin{cases} S(C)_{\text{lin}}(\rho) & \rho \in L \\ S(C')_{\text{lin}}(\rho) & \text{otherwise} \end{cases} \\
S(\text{Filter}(C, L))_{\text{lin}}(\rho) &= \begin{cases} S(C)_{\text{lin}}(\rho) & \rho \in L \\ 0 & \text{otherwise} \end{cases} \\
S(\text{Assign}(C, \rho' : \tau))_{\text{lin}}(\rho) &= S(C)_{\text{lin}}(\rho)
\end{aligned}$$

(b) Linearities

$$S(C)_{\text{lin}}(\rho) = \omega \implies S(C)(\rho) \leq S(\text{Assign}(C, \rho : \tau))(\rho) \text{ for all stores } \text{Assign}(C, \rho : \tau)$$

(c) Weak updates

$$\begin{aligned}
&\frac{x \in \text{dom}(\Gamma)}{\Gamma, C \vdash x : \Gamma(x), C} \text{ (Var)} \\
&\frac{\kappa \text{ fresh}}{\Gamma, C \vdash n : \kappa \text{ int}, C} \text{ (Int)} \\
&\frac{\Gamma, C \vdash e : \tau, C' \quad \tau \leq C'(\rho) \quad \kappa \text{ fresh}}{\Gamma, C \vdash \text{ref}^\rho e : \kappa \text{ ref}(\rho), \text{Alloc}(C', \rho)} \text{ (Ref)} \\
&\frac{\Gamma, C \vdash e : Q \text{ ref}(\rho), C'}{\Gamma, C \vdash !e : C'(\rho), C'} \text{ (Deref)} \\
&\frac{\Gamma, C \vdash e_1 : Q \text{ ref}(\rho), C' \quad \Gamma, C' \vdash e_2 : \tau, C''}{\Gamma, C \vdash e_1 := e_2 : \tau, \text{Assign}(C'', \rho : \tau)} \text{ (Assign)} \\
&\frac{\tau = \text{sp}(t) \quad \varepsilon, \varepsilon', \kappa \text{ fresh} \quad \Gamma[x \mapsto \tau], \varepsilon \vdash e : \tau', C' \quad C' \leq \varepsilon'}{\Gamma, C \vdash \lambda^L x.t.e : \kappa(\varepsilon, \tau) \longrightarrow^L (\varepsilon', \tau'), C} \text{ (Lam)} \\
&\frac{\Gamma, C \vdash e_1 : Q(\varepsilon, \tau) \longrightarrow^L (\varepsilon', \tau'), C' \quad \Gamma, C' \vdash e_2 : \tau_2, C'' \quad \tau_2 \leq \tau \quad \text{Filter}(C'', L) \leq \varepsilon}{\Gamma, C \vdash e_1 e_2 : \tau', \text{Merge}(\varepsilon', C'', L)} \text{ (App)} \\
&\frac{\Gamma, C \vdash e : Q' \sigma, C' \quad Q' \leq Q}{\Gamma, C \vdash \text{assert}(e, Q) : Q \sigma, C'} \text{ (Assert)} \\
&\frac{\Gamma, C \vdash e : Q' \sigma, C' \quad Q' \leq Q}{\Gamma, C \vdash \text{check}(e, Q) : Q' \sigma, C'} \text{ (Check)}
\end{aligned}$$

Figure 5: Constraint generation rules

Figure 4: Extending a solution to constructed stores

C , evaluating e yields a result of type τ and a new store C' . We write $C(\rho)$ for the type associated with ρ in store C ; we discuss the computation of $C(\rho)$ in Section 3.4. We use the function $\text{sp}(t)$ to decorate a standard type t with fresh qualifier and store variables:

$$\begin{aligned}
\text{sp}(\alpha) &= \kappa \alpha & \kappa \text{ fresh} \\
\text{sp}(\text{int}) &= \kappa \text{ int} & \kappa \text{ fresh} \\
\text{sp}(\text{ref}(\rho)) &= \kappa \text{ ref}(\rho) & \kappa \text{ fresh} \\
\text{sp}(t \longrightarrow^L t') &= \kappa(\varepsilon, \text{sp}(t)) \longrightarrow^L (\varepsilon', \text{sp}(t')) & \kappa, \varepsilon, \varepsilon' \text{ fresh}
\end{aligned}$$

We briefly discuss the rules in Figure 5:

- (Var) and (Int) are standard. For (Int), we pick a fresh qualifier variable κ to annotate n 's type.
- (Ref) adds a location ρ to the store C' , yielding the store $\text{Alloc}(C', \rho)$. The type τ of e is constrained to be compatible with ρ 's type in C' .²
- (Deref) looks up the type of e 's location ρ in the current store C' . In this rule, any qualifier may appear on e 's type; qualifiers are checked only by (Check), see below.
- (Assign) produces a new store representing the assignment of type τ to location ρ .
- (Lam) type checks function body e in fresh initial store ε and with parameter x bound to a type with fresh qualifier variables.

²An alternative formulation is to track the type τ of e as part of the constructed store $\text{Alloc}(\cdot, \cdot)$, and only constrain τ to be compatible with $C'(\rho)$ if after the allocation ρ is non-linear.

- (App) constrains $\tau_2 \leq \tau$ to ensure that e_2 's type is compatible with e_1 's argument type. The constraint $\text{Filter}(C'', L) \leq \varepsilon$ ensures that the current state of the locations that e_1 uses, which are captured by its effect set L , is compatible with the state function e_1 expects. The final store $\text{Merge}(\varepsilon', C'', L)$ joins the store C'' before the function call with the result store ε' of the function. Intuitively, this rule gives us some low-cost polymorphism, in which functions do not act as join points for locations they do not use.
- (Assert) adds a qualifier annotation to the program, and (Check) checks that the inferred top-level qualifier Q' of e is compatible with the expected qualifier Q .

Figure 6 shows the stores and store constraints generated for our example program. We have slightly simplified the graph for clarity. Here ε is f 's initial store and ε' is f 's final store. We use undirected edges for store constructors and a directed edge from C_1 to C_2 for the constraint $C_1 \leq C_2$.

We step through constraint generation. We model the allocation of ρ_x with the store $\text{Alloc}(\varepsilon, \rho_x)$. Location ρ_x is initialized to 0, which is given the type $\kappa_0 \text{ int}$ for fresh qualifier variable κ_0 . (Ref) generates the constraint $\kappa_0 \text{ int} \leq \varepsilon(\rho_x)$ to require that the type of 0 be compatible with $\varepsilon(\rho_x)$. We model the allocation and initialization of ρ_y and ρ_z similarly. Then we construct three *Assign* stores to represent the assignment statements. We give 3 and 4 the types $\kappa_3 \text{ int}$ and $\kappa_4 \text{ int}$, respectively, where κ_3 and κ_4 are fresh qualifier variables.

For the recursive call to f , we construct a *Filter* and add a constraint on ε . The *Merge* store represents the state when the recursive call to f returns. We join the two branches of

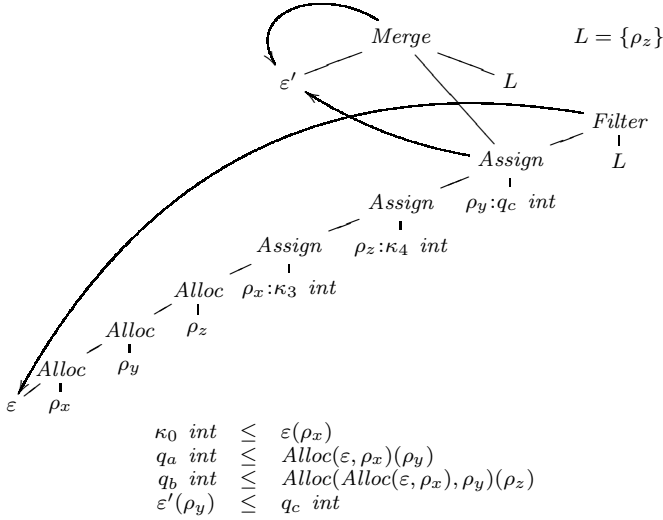


Figure 6: Store constraints for example

the conditional by making edges to ε' . Notice the cycle, due to recursion, in which state from ε' can flow to the *Merge*, which in turn can flow to ε' . Finally, the qualifier check requires that $\varepsilon'(\rho_y)$ has qualifier q_c .

3.4 Flow-Sensitive Constraint Resolution

The rules of Figure 5 generate three kinds of constraints: qualifier constraints $Q \leq Q'$, subtyping constraints $\tau \leq \tau'$, and store constraints $C \leq \varepsilon$ (the right-hand side of a store constraint is always a store variable). A set of m type and qualifier constraints can be solved in $O(m)$ time using well-known techniques [16, 23], so in this section we focus on computing a solution S to a set of store constraints.

Our analysis is most precise if as few locations as possible are non-linear. Recall that linearities naturally form a partial order $0 < 1 < \omega$. Thus, given a set of constructed stores and store constraints, we perform a least fixpoint computation to determine $S(C)_{lin}(\rho)$. We initially assume that in every store, location ρ has linearity 0. Then we exhaustively apply the rules in Figure 4(b) and the rule $S(\varepsilon)_{lin}(\rho) = (\max_{\{C|C \leq \varepsilon\}} S(C)_{lin}(\rho))$ until we reach a fixpoint. This last rule is derived from Figure 3.

In our implementation, we compute $S(C)_{lin}(\rho)$ in a single pass over the store constraints using Tarjan's strongly-connected components algorithm to find cycles in the store constraint graph. For each such cycle containing more than one allocation of the same location ρ we set the linearity of ρ to ω in all stores on the cycle.

Given this algorithm to compute $S(C)_{lin}(\rho)$, in principle we can then solve the implied typing constraints using the following simple procedure. For each store variable ε , initialize $S(\varepsilon)$ to a map

$$\{\rho_1:sp(C_I(\rho_1)), \dots, \rho_n:sp(C_I(\rho_n))\}$$

and for each store $Assign(C, \rho : \tau)$ initialize $S(Assign(C, \rho : \tau))(\rho)$ to $sp(C_I(\rho))$, thereby assigning fresh qualifiers to the type of every location at every program point. Replace uses of $C(\rho)$ in Figure 5 with $S(C)(\rho)$, using the logic in Figure 4(a).

Apply the following two closure rules until no more con-

straints are generated:

$$\begin{aligned} C \leq \varepsilon &\implies S(C)(\rho) \leq S(\varepsilon)(\rho) && \text{for all } \rho \\ S(C)_{lin}(\rho) = \omega &\implies S(C)(\rho) \leq S(Assign(C, \rho : \tau))(\rho) && \text{for all stores } Assign(C, \rho : \tau) \end{aligned}$$

Given a program of size n , in the worst case this naive algorithm requires at least n^2 space and time to build $S(\cdot)$ and generate the necessary type constraints. This cost is too high for all but small examples. We reduce this cost in practice by taking advantage of several observations.

Many locations are flow-insensitive. If a location ρ never appears on the left-hand side of an assignment, then ρ 's type cannot change. Thus we can give ρ one global type instead of one type per program point. In imperative languages such as C, C++, and Java, function parameters are a major source of flow-insensitive locations. In these languages, because parameters are *l*-values, they have an associated memory location that is initialized but then often never subsequently changed.

Adding extra store variables trades space for time. To compute $S(C)(\rho)$ for a constructed store C , we must deconstruct C recursively until we reach a variable store or an assignment to ρ (see Figure 4(a)). Because we represent the effect constraints compactly (in linear space), deconstructing $Filter(\cdot, L)$ or $Merge(\cdot, \cdot, L)$ may require a potentially linear time computation to check whether $\rho \in L$. We recover efficient lookups by replacing C with a fresh store variable ε and adding the constraint $C \leq \varepsilon$. Then rather than computing $S(C)(\rho)$ we compute $S(\varepsilon)(\rho)$, which requires only a map lookup. Of course, we must use space to store ρ in $S(\varepsilon)$. However, as shown below, we often can avoid this cost completely. We apply this transformation to each store $Merge(C, C', L)$ constructed during constraint inference.

Not every store needs every location. Rather than assuming $S(\varepsilon)$ contains all locations, we add needed locations lazily. We add a location ρ to $S(\varepsilon)$ the first time the analysis requests $\varepsilon(\rho)$ and whenever there is a constraint $C \leq \varepsilon$ or $\varepsilon \leq C$ such that $\rho \in S(C)$. Stores constructed with *Filter* and *Merge* will tend to stop propagation of locations, saving space (e.g., if $Filter(C, L) \leq \varepsilon$, $\rho \in S(\varepsilon)$, but $\rho \notin L$, then we do not propagate ρ to C).

We can extend this idea further. For each qualifier variable κ , inference maintains a set of possible qualifier constants that are valid solutions for κ . If that set contains every constant qualifier, then κ is *uninteresting* (i.e., κ is constrained only by other qualifier variables), otherwise κ is *interesting*. A type τ is interesting if any qualifier in τ is interesting, otherwise τ is uninteresting. We then modify the closure rules as follows:

$$\begin{aligned} C \leq \varepsilon &\implies S(C)(\rho) \leq S(\varepsilon)(\rho) \\ &\text{for all } \rho \in S(C) \text{ or } S(\varepsilon) \text{ s.t.} \\ &S(C)(\rho) \text{ or } S(\varepsilon)(\rho) \text{ interesting} \end{aligned}$$

$$\begin{aligned} S(C)_{lin}(\rho) = \omega &\implies S(C)(\rho) \leq S(Assign(C, \rho : \tau))(\rho) \\ &\text{for all } Assign(C, \rho : \tau) \text{ s.t. } S(C)(\rho) \text{ or} \\ &S(Assign(C, \rho : \tau))(\rho) \text{ interesting} \end{aligned}$$

In this way, if a location ρ is bound to an uninteresting type, then we need not propagate ρ through the constraint graph.

Figure 7 gives an algorithm for lazy location propagation. We associate a mark with each ρ in each $S(\varepsilon)$ and with ρ in $Assign(C, \rho : \tau)$. Initially this mark is not set, indicating that location ρ is bound to an uninteresting type.

If a qualifier variable κ appears in $S(\varepsilon)(\rho)$, we associate the pair (ρ, ε) with κ , and similarly for $Assign$ stores. If during constraint resolution the set of possible solutions of κ changes, we call $PROPAGATE(\rho, C)$ to propagate ρ , and in turn κ , through the store constraint graph.

If $PROPAGATE(\rho, C)$ is called and ρ is already marked in C , we do nothing. Otherwise, $BACK-PROP()$ and $FORWARD-PROP()$ make appropriate constraints between $S(C)(\rho)$ and $S(C')(\rho)$ for every store C' reachable from C . This step may add ρ to C' if C' is a store variable, and the type constraints that $BACK-PROP()$ and $FORWARD-PROP()$ generate may trigger subsequent calls to $PROPAGATE()$.

Consider again our running example. Figure 8 shows how locations and qualifiers propagate through the store constraint graph. Dotted edges in this graph indicate inferred constraints (discussed below). For clarity we have omitted the $Alloc$ edges (summarized with a dashed line) and the base types.

The four type constraints in Figure 6 are shown as directed edges in Figure 8. For example, the constraint $\kappa_0 \text{ int} \leq \varepsilon(\rho_x)$ reduces to the constraint $\kappa_0 \leq \kappa_x$, which is a directed edge $\kappa_0 \rightarrow \kappa_x$. Adding this constraint does not cause any propagation; this constraint is among variables. Notice that the assignment of type $\kappa_3 \text{ int}$ to ρ_x also does not cause any propagation.

The constraint $q_a \text{ int} \leq Alloc(\varepsilon, \rho_x)(\rho_y)$ reduces to $q_a \text{ int} \leq \varepsilon(\rho_y)$, which reduces to $q_a \leq \kappa_y$. This constraint does trigger propagation. $PROPAGATE(\rho_y, \varepsilon)$ first pushes ρ_y backward to the $Filter$ store. But since $\rho_y \notin L$, propagation stops. Next we push ρ_y forward through the graph and stop when we reach the store $Assign(\cdot, \rho_y : q_c \text{ int})$; forward propagation assumes that this is a strong update.

Since $Assign(\cdot, \rho_y : q_c \text{ int})$ contains an interesting type, ρ_y is propagated from this store forward through the graph. On one path, propagation stops at the $Filter$. The other paths yield a constraint $q_c \leq \kappa'_y$. Notice that the constraint $\kappa'_y \leq q_c$ remains satisfiable.

The constraint $q_b \leq \kappa_z$ triggers a propagation step as before. However, this time $\kappa_z \in L$, and during backward propagation when we reach $Filter$ we must continue. Eventually we reach $Assign(\cdot, \rho_z : \kappa_4 \text{ int})$ and add the constraint $\kappa_4 \leq \kappa_z$. This in turn triggers propagation from $Assign(\cdot, \rho_z : \kappa_4 \text{ int})$. This propagation step reaches ε' , adds ρ_z to $S(\varepsilon')$, and generates the constraint $\kappa_4 \leq \kappa'_z$.

Finally, we determine that in the $Assign$ stores ρ_x and ρ_y are linear and ρ_z is non-linear. (The linearity computation uses the $Alloc(\cdot, \cdot)$ stores, which are not shown.) Thus the update to ρ_z is a weak update, which yields a constraint $\kappa_z \leq \kappa_4$.

This example illustrates three kinds of propagation. The location ρ_x is never interesting, so it is not propagated through the graph. The location ρ_y is propagated, but propagation stops at the strong update to ρ_y and also at the $Filter$, because the (Down) rule in Figure 1 was able to prove that ρ_y is purely local to f . The location ρ_z , on the other hand, is not purely local to f , and thus all instances of ρ_z are conflated, and ρ_z admits only weak updates.

```

PROPAGATE( $\rho, \varepsilon$ ) =
  case C of
     $\varepsilon$ :
      add  $\rho : sp(C_I(\rho))$  to  $S(\varepsilon)$  if not already in  $S(\varepsilon)$ 
      if  $\rho$  is not marked in  $\varepsilon$ 
        mark  $\rho$  in  $S(\varepsilon)$ 
        FORWARD-PROP( $C, \rho, S(\varepsilon)(\rho)$ )
        for each  $C'$  such that  $C' \leq \varepsilon$ 
          BACK-PROP( $C', \rho, S(\varepsilon)(\rho)$ )
      Assign( $C', \rho : \tau$ ):
        if  $\rho$  is not marked in Assign( $C', \rho : \tau$ )
          mark  $\rho$  in Assign( $C', \rho : \tau$ )
          FORWARD-PROP( $C, \rho, \tau$ )

BACK-PROP( $C, \rho, \tau$ ) =
  case C of
     $\varepsilon$ :
      add  $\rho : sp(C_I(\rho))$  to  $S(\varepsilon)$  if not already in  $S(\varepsilon)$ 
       $S(\varepsilon)(\rho) \leq \tau$ 
      Alloc( $C', \rho'$ ):
        BACK-PROP( $C', \rho, \tau$ )
      Merge( $C', C'', L$ ):
        if  $\rho \in L$ 
          then BACK-PROP( $C', \rho, \tau$ )
          else BACK-PROP( $C'', \rho, \tau$ )
      Filter( $C', L$ ):
        if  $\rho \in L$ 
          then BACK-PROP( $C', \rho, \tau$ )
      Assign( $C', \rho' : \tau'$ ):
        if  $\rho = \rho'$ 
          then  $\tau' \leq \tau$ 
          else BACK-PROP( $C', \rho, \tau$ )

FORWARD-PROP( $C, \rho, \tau$ ) =
  for each  $\varepsilon$  such that  $C \leq \varepsilon$ 
    add  $\rho : sp(C_I(\rho))$  to  $S(\varepsilon)$  if not already in  $S(\varepsilon)$ 
     $\tau \leq S(\varepsilon)(\rho)$ 
  for each  $C'$  such that  $C'$  is constructed from C
    case C' of
      Alloc( $C', \rho'$ ):
        FORWARD-PROP( $C', \rho, \tau$ )
      Merge( $C_1, C_2, L$ ):
        if  $\rho \in L$  and  $C = C_1$ 
          then FORWARD-PROP( $C', \rho, \tau$ )
        if  $\rho \notin L$  and  $C = C_2$ 
          then FORWARD-PROP( $C', \rho, \tau$ )
      Filter( $C, L$ ):
        if  $\rho \in L$ 
          then FORWARD-PROP( $C', \rho, \tau$ )
      Assign( $C, \rho' : \tau'$ ):
        if  $\rho \neq \rho'$ 
          then FORWARD-PROP( $C', \rho, \tau$ )

```

Figure 7: Lazy location constraint propagation

4. RESTRICT

As mentioned in the introduction, type inference may fail because a location on which a strong update is needed may be non-linear. In practice a major source of non-linear locations is data structures. For example, given a linked list l , our alias analysis often cannot distinguish $l \rightarrow \text{lock}$ from $l \rightarrow \text{next} \rightarrow \text{lock}$, hence both will likely be non-linear.

Our solution to this problem is to add a new form

```
restrict  $x = e_1$  in  $e_2$ 
```

to the language. Intuitively, this declares that of all aliases of e_1 , only x and copies derived from x will be used within

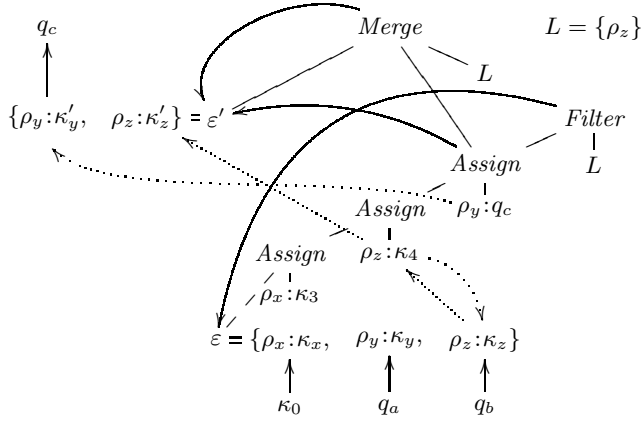


Figure 8: Constraint propagation

e_2 . For example, consider

```
restrict x = y in {
  x := ...; /* valid */
  y := ...; /* invalid */
}
```

The first assignment through x is valid, but the assignment through y is forbidden by `restrict`.

We check `restrict` using the following type rule, which is integrated into the first inference pass of Figure 1:

$$\frac{\Gamma \vdash e_1 \Rightarrow e'_1 : t_1; L_1 \quad t_1 = \text{ref}(\rho) \quad \rho, \rho' \text{ fresh} \quad C_I(\rho') = C_I(\rho) \quad \Gamma[x \mapsto \text{ref}(\rho')] \vdash e_2 \Rightarrow e'_2 : t_2; L_2 \quad \rho \notin L_2 \quad \rho' \notin \text{locs}(\Gamma) \cup \text{locs}(C_I(\rho)) \cup \text{locs}(t_2)}{\Gamma \vdash \text{restrict } x = e_1 \text{ in } e_2 \Rightarrow \text{restrict}^{\rho'} x = e'_1 \text{ in } e'_2 : t_2; L_1 \cup L_2 \cup \{\rho\}} \text{ (Restrict)}$$

Here we bind x to a type with a fresh abstract location ρ' to distinguish dereferences of x from dereferences of other aliases of e_1 . The constraint $\rho \notin L_2$ forbids location ρ from being dereferenced in e_2 ; notice dereferences of ρ' within e_2 are allowed. We require that ρ' not escape the scope of e_2 with $\rho' \notin \text{locs}(\Gamma) \cup \text{locs}(C_I(\rho)) \cup \text{locs}(t_2)$, and we also add ρ to the effect set. We translate `restrict` into the target language by annotating it with the location ρ' that x is bound to. A full discussion of `restrict`, including a soundness proof, can be found in a technical report [15].

We use `restrict` to locally recover strong updates. The key observation is that the location ρ of e_1 and the location ρ' of x can be different. Thus even if the linearity of ρ is ω , the linearity of ρ' can be 1. Therefore within the body of e_2 we may be able to perform strong updates of ρ' . When the scope of `restrict` ends, we may need to do a weak update from ρ' to ρ .

For example, suppose that we wish to type check a state change of some lock deep within a data structure, and the location of the lock is non-linear. The following is not atypical of Linux kernel code:

```
spin_lock(&a->b[c].d->lock); /* invalid; */
... /* non-linear loc */
spin_unlock(&a->b[c].d->lock);
```

Assuming the type system determines that the `...` above contains no accesses to aliases of the lock and does not alias

the lock to a non-linear location, we can modify the code to type check as follows:

```
restrict lock = &a->b[c].d->lock in {
  spin_lock(lock); /* valid */
  ...
  spin_unlock(lock);
}
```

In our flow-sensitive step, we use the following inference rule for `restrict`:

$$\frac{\Gamma, C \vdash e_1 : Q \text{ ref}(\rho), C' \quad C'' = \text{Alloc}(C', \rho') \quad C'(\rho) \leq C''(\rho') \quad \Gamma[x \mapsto \text{ref}(\rho')], C'' \vdash e_2 : \tau_2, C'''}{\Gamma, C \vdash \text{restrict}^{\rho'} x = e_1 \text{ in } e_2 : \tau_2, \text{Assign}(C''', \rho : C'''(\rho'))} \text{ (Restrict)}$$

In this rule, we infer a type for e_1 , which is a pointer to some location ρ . Then we create a new store C'' in which the location ρ' of x is both allocated and initialized to $C'(\rho)$. In C'' , and with x added to the type environment, we evaluate e_2 . Finally, the result store is the store C''' with a potentially weak update assigning the contents of ρ' to ρ .

5. EXPERIMENTS

To test our ideas in practice we have built a tool CQUAL that implements our inference algorithm. To use CQUAL, programmers annotate their C programs with type qualifiers, which are added to the C syntax in the same way as `const` [16]. The tool CQUAL can analyze a single file or a whole program. As is standard in type-based analysis, when analyzing a single file, the programmer supplies type signatures for any external functions or variables.

We have used CQUAL to check two program properties: locking in the 2.4.9 Linux kernel device drivers and uses of the C stream library. Our implementation is sound up to the unsafe features of C: type casts, variable-argument functions, and ill-defined pointer arithmetic. We currently make no attempt to track the effect of any of these features on aliasing, except for the special case of type casting the result of `malloc`-like functions. In combination with a system for enforcing memory safety, such as CCured [21], our implementation would be sound.

In our implementation, we do not allow strong updates on locations containing functions. This improves efficiency because we never need to recompute $S(C)_{\text{lin}}(\rho)$ —weak updates will not add constraints between stores. Additionally, observe that allocations affect linearities but not types, and reads and writes affect types but not linearities. Thus in our implementation we also improve the precision of the analysis by distinguishing read, write, and allocation effects. We omit details due to space constraints.

The analysis results are presented to the user with an `emacs`-based user interface. The source code is colored according to the inferred qualifiers. Type errors are hyper-linked to the source line at which the error first occurred, and the user can click on qualifiers to view a path through the constraint graph that shows why a type error was detected. We have found the ability to visualize constraint solutions in terms of the original source syntax not just useful, but essential, to understanding the results of inference. More detail on the ideas in the user interface can be found in [24].

5.1 Linux Kernel Locking

The Linux kernel includes two primitive locking functions, which are used extensively by device drivers:

```
void spin_lock(spinlock_t *lock);
void spin_unlock(spinlock_t *lock);
```

We use three qualifiers `locked`, `unlocked`, and \top (unknown) to check locking behavior. The subtyping relation is `locked` < \top and `unlocked` < \top . We assign `spin_lock` the type

$$(C, \text{ref}(\rho)) \longrightarrow^{\{\rho\}} (\text{Assign}(C, \rho : \text{locked spinlock_t}), \text{void})$$

where

$$C(\rho) \leq \text{unlocked spinlock_t}$$

We omit the function qualifier since it is irrelevant. The type of `spin_lock` requires that the lock passed as the argument be `unlocked` (see the where clause) and changes it to `locked` upon returning. The signature for `spin_unlock` is the same with `locked` and `unlocked` exchanged.

In practice we give `spin_lock` this type signature by supplying CQUAL with the following definition:

```
void spin_lock($unlocked spinlock_t *lock) {
  change_type(*lock, $locked spinlock_t);
}
```

Here `change_type(x, t)` is just like the assignment

$$x = \langle \text{something of type } t \rangle;$$

except that rather than give an explicit right-hand side we just give the type of the right-hand side. In this case the programmer needs to supply the body of `spin_lock` because it is inline assembly code.

Since our implementation currently lacks parametric polymorphism, we inline calls to `spin_lock` and `spin_unlock`.

Using these type signatures we can check for three kinds of errors: deadlocks from acquiring a lock already held by the same thread, attempts to release a lock already released by the same thread, and attempting to acquire or release a lock in an unknown (\top) state.

We analyzed 513 whole device driver modules (a whole module includes all the files that make up a single driver). A module must meet a well-specified kernel interface, which we model with a `main` function that non-deterministically calls all possible driver functions registered with the kernel.

We also separately analyzed each of the 892 driver files making up the whole modules. In these experiments we removed the \top qualifier so that `locked` and `unlocked` are incomparable, and we made optimistic assumptions about the environment in which each file is invoked.

We examined the results for 64 of the 513 whole device driver modules and for all of the 892 separately analyzed driver files. We found 14 apparently new locking bugs, including one which spans multiple files. In five of the apparent bugs a function tries to acquire a lock already held by a function above it in the call chain, leading to a deadlock. For example, the `emu10k1` module contains a deadlock (we omit the `void` return types):

```
emu10k1_mute_irqhandler(struct emu10k1_card *card) {
  struct patch_manager *mgr = &card->mgr;
  ... spin_lock_irqsave(&mgr->lock, flags);
  emu10k1_set_oss_vol(card, ...); ...
}
emu10k1_set_oss_vol(struct emu10k1_card *card, ...) {
  ... emu10k1_set_volume_gpr(card, ...); ...
}
```

```
emu10k1_set_volume_gpr(struct emu10k1_card *card, ...) {
  struct patch_manager *mgr = &card->mgr;
  ... spin_lock_irqsave(&mgr->lock, flags); ...
}
```

Note that detecting this error requires interprocedural analysis.

One of our goals is to understand how often, and why, our system fails to type check real programs. We have categorized every type error in the separate file analysis of the 892 driver files. In this experiment, of the 52 files that fail to type check, 11 files have locking bugs (sometimes more than one) and the remaining 41 files have type errors. Half of these type errors are due to incorrect assumptions about the interface for functions; these type errors are eliminated by moving to whole module analysis. The remaining type errors fall into two main categories.

In many cases the problem is that our alias analysis is not strong enough to type check the program. Another common class of type errors arises when locks are conditionally acquired and released. In this case, a lock is acquired if a predicate P is true. Before the lock is released, P is tested again to check whether the lock is held. Our system is not path sensitive, and our tool signals a type error at the point where the path on which the lock is acquired joins with the path on which the lock is not acquired (since we did not use \top in these single file experiments—in the whole module analysis, this error is detected later on, when there is an attempt to acquire or release the lock in the \top state). Most of these examples could be rewritten with little effort to pass our type system. In our opinion, this would usually make the code clearer and safer—the duplication of the test on P invites new bugs when the program is modified.

Even after further improvements, we expect some dynamically correct programs will not type check. As future work, we propose the following solution. The qualifier \top represents an unknown state. We can use the information in the constraints to automatically insert coercions to and from \top where needed. During execution these coercions perform runtime tests to verify locks are in the correct state. Thus, our approach can introduce dynamic type checking in situations where we cannot prove safety statically.

Of the 513 whole modules, 196 contain type errors, many of which are duplicates from shared code. We examined 64 of the type error-containing modules and discovered that a major source of type errors is when there are multiple aliases of a location, but only one alias is actually used in the code of interest. Not surprisingly, larger programs, such as whole modules, have more problems with spurious aliasing than the optimistic single-file analysis. We added `restrict` annotations by hand to the 64 modules we looked at, including the `emu10k1` module, which yielded the largest number of such false positives. Using `restrict`, we eliminated all of the false positives in these modules that occurred because non-linear locations could not be strongly updated. This supports our belief that `restrict` is the right tool for dealing with (necessarily) conservative alias analysis. Currently adding `restrict` by hand is burdensome, requiring a relatively large number of annotations. We leave the problem of automatically inferring `restrict` annotations as future work.

5.2 C Stream Library

As mentioned in the introduction, the C stream library

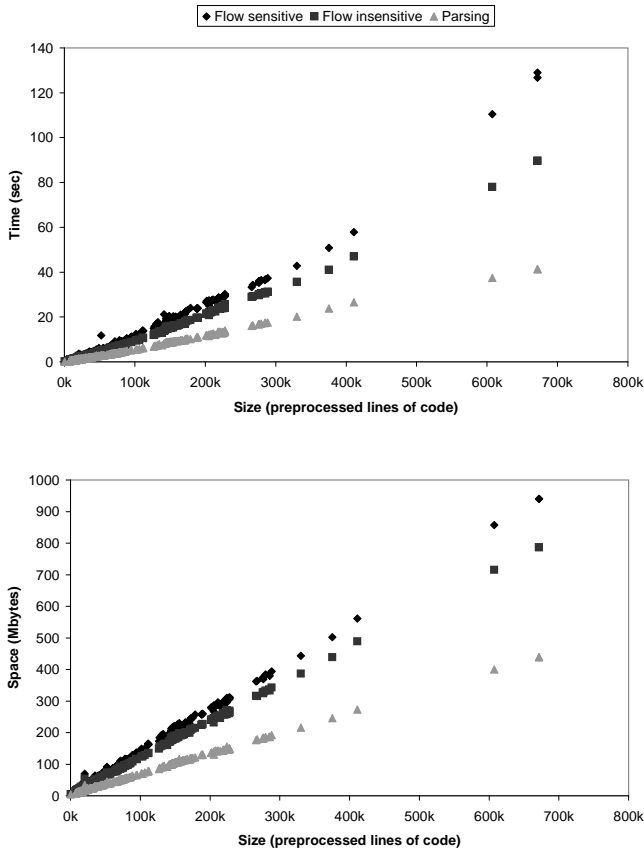


Figure 9: Resource usage for whole module analysis

interface contains certain sequencing constraints. For example, a file must be opened for reading before being read. A special property of the C stream library is that the result of `fopen` must be tested against `NULL` before being used, because `fopen` may or may not succeed. The class of C stream library usage errors our tool can detect includes files used without having been opened and checked against `NULL`, files opened and then accessed in an incompatible mode, and files accessed after being closed. We omit the details due to space constraints.

We tried our tool on two application programs, `man-1.5h1` and `sendmail-8.11.6`. We were primarily interested in the performance of our tool on a more complex application (see below), as we did not expect to find any latent stream library usage bugs in such mature programs. However, we did find one minor bug in `sendmail`, in which an opened log file is never closed in some circumstances.

5.3 Precision and Efficiency

The algorithm described in Section 3.4 is carefully designed to limit resource usage. Figure 9 shows time and space usage of whole module analysis versus preprocessed lines of code for 513 Linux kernel modules. All experiments were done on a dual processor 550 MHz Pentium III with 2GB of memory running RedHat 6.2.

We divide the resource usage into C parsing and type checking, flow-insensitive analysis, and flow-sensitive analysis. Flow-insensitive analysis consists of the alias and effect inference of Figure 1 together with flow-insensitive qualifier

inference [16]. Flow-sensitive analysis consists of the constraint generation and resolution described in Sections 3.3-3.4, including the linearity computation. In the graphs, the reported time and space for each phase includes the time and space for the previous phases.

The graphs show that the space overhead of flow-sensitive analysis is relatively small and appears to scale well to large modules. For all modules the space usage for the flow-sensitive analysis is within 31% of the space usage for the flow-insensitive analysis. The running time of the analysis is more variable, but the absolute running times are within a factor of 1.3 of the flow-insensitive running times.

The analysis of `sendmail-8.11.6`, with 175,193 preprocessed source lines, took 28.8 seconds and 264MB; `man-1.5h1`, with 16,411 preprocessed source lines, took 1.85 seconds and 32MB. These results suggest that our algorithm also behaves efficiently when checking C stream library usage.

6. CONCLUSION

We have presented a system for extending standard type systems with flow-sensitive type qualifiers. We have given a lazy constraint resolution algorithm to infer type qualifier annotations and have shown that our analysis is effective in practice by finding a number of new locking bugs in the Linux kernel.

7. REFERENCES

- [1] A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison Wesley, 1988.
- [2] A. Aiken, M. Fähndrich, and R. Levien. Better Static Memory Management: Improving Region-Based Analysis of Higher-Order Languages. In *Proceedings of the 1995 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 174–185, La Jolla, California, June 1995.
- [3] R. Altucher and W. Landi. An Extended Form of Must Alias Analysis for Dynamic Allocation. In *Proceedings of the 22nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 74–84, San Francisco, California, Jan. 1995.
- [4] D. R. Chase, M. Wegman, and F. K. Zadeck. Analysis of Pointers and Structures. In *Proceedings of the 1990 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 296–310, White Plains, New York, June 1990.
- [5] K. Cray, D. Walker, and G. Morrisett. Typed Memory Management in a Calculus of Capabilities. In *Proceedings of the 26th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 262–275, San Antonio, Texas, Jan. 1999.
- [6] M. Das, S. Lerner, and M. Seigle. ESP: Path-Sensitive Program Verification in Polynomial Time. In *Proceedings of the 2002 ACM SIGPLAN Conference on Programming Language Design and Implementation*, Berlin, Germany, June 2002. To appear.
- [7] R. DeLine and M. Fähndrich. Enforcing High-Level Protocols in Low-Level Software. In *Proceedings of the 2001 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 59–69, Snowbird, Utah, June 2001.
- [8] M. Emami, R. Ghiya, and L. J. Hendren. Context-Sensitive Interprocedural Points-to Analysis in the Presence of Function Pointers. In *Proceedings of the 1994 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 242–256, Orlando, Florida, June 1994.
- [9] D. Engler, B. Chelf, A. Chou, and S. Hallem. Checking System Rules Using System-Specific, Programmer-Written

- Compiler Extensions. In *Fourth symposium on Operating System Design and Implementation*, San Diego, California, Oct. 2000.
- [10] D. Engler, D. Y. Chen, S. Hallem, A. Chou, and B. Chelf. Bugs as Deviant Behavior: A General Approach to Inferring Errors in Systems Code. In *Proceedings of the 18th ACM Symposium on Operating Systems Principles*, Banff, Canada, Oct. 2001.
- [11] D. Evans. Static Detection of Dynamic Memory Errors. In *Proceedings of the 1996 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 44–53, Philadelphia, Pennsylvania, May 1996.
- [12] M. Fähndrich and R. DeLine. Adoption and Focus: Practical Linear Types for Imperative Programming. In *Proceedings of the 2002 ACM SIGPLAN Conference on Programming Language Design and Implementation*, Berlin, Germany, June 2002. To appear.
- [13] C. Flanagan and S. N. Freund. Type-Based Race Detection for Java. In *Proceedings of the 2000 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 219–232, Vancouver B.C., Canada, June 2000.
- [14] C. Flanagan, K. R. M. Leino, M. Lillibridge, G. Nelson, J. B. Saxe, and R. Stata. Extended Static Checking for Java. In *Proceedings of the 2002 ACM SIGPLAN Conference on Programming Language Design and Implementation*, Berlin, Germany, June 2002. To appear.
- [15] J. S. Foster and A. Aiken. Checking Programmer-Specified Non-Aliasing. Technical Report UCB//CSD-01-1160, University of California, Berkeley, Oct. 2001.
- [16] J. S. Foster, M. Fähndrich, and A. Aiken. A Theory of Type Qualifiers. In *Proceedings of the 1999 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 192–203, Atlanta, Georgia, May 1999.
- [17] D. Grossman, G. Morrisett, Y. Wang, T. Jim, M. Hicks, and J. Cheney. Cyclone user's manual. Technical Report 2001-1855, Department of Computer Science, Cornell University, Nov. 2001. Current version at <http://www.cs.cornell.edu/projects/cyclone>.
- [18] A. Igarashi and N. Kobayashi. Resource Usage Analysis. In *Proceedings of the 29th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 331–342, Portland, Oregon, Jan. 2002.
- [19] S. Jagannathan, P. Thiemann, S. Weeks, and A. Wright. Single and loving it: Must-alias analysis for higher-order languages. In *Proceedings of the 25th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 329–341, San Diego, California, Jan. 1998.
- [20] J. M. Lucassen and D. K. Gifford. Polymorphic Effect Systems. In *Proceedings of the 15th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 47–57, San Diego, California, Jan. 1988.
- [21] G. Necula, S. McPeak, and W. Weimer. CCured: Type-Safe Retrofitting of Legacy Code. In *Proceedings of the 29th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 128–139, Portland, Oregon, Jan. 2002.
- [22] R. O'Callahan. A Simple, Comprehensive Type System for Java Bytecode Subroutines. In *Proceedings of the 26th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 70–78, San Antonio, Texas, Jan. 1999.
- [23] J. Rehof and T. Æ. Mogensen. Tractable Constraints in Finite Semilattices. In R. Cousot and D. A. Schmidt, editors, *Static Analysis, Third International Symposium*, volume 1145 of *Lecture Notes in Computer Science*, pages 285–300, Aachen, Germany, Sept. 1996. Springer-Verlag.
- [24] U. Shankar, K. Talwar, J. S. Foster, and D. Wagner. Detecting Format String Vulnerabilities with Type Qualifiers. In *Proceedings of the 10th Usenix Security Symposium*, Washington, D.C., Aug. 2001.
- [25] F. Smith, D. Walker, and G. Morrisett. Alias Types. In G. Smolka, editor, *9th European Symposium on Programming*, volume 1782 of *Lecture Notes in Computer Science*, pages 366–381, Berlin, Germany, 2000. Springer-Verlag.
- [26] R. Stata and M. Abadi. A Type System for Java Bytecode Subroutines. In *Proceedings of the 25th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 149–160, San Diego, California, Jan. 1998.
- [27] R. E. Strom and S. Yemini. Typestate: A Programming Language Concept for Enhancing Software Reliability. *IEEE Transactions on Software Engineering*, 12(1):157–171, Jan. 1986.
- [28] M. Tofte and J.-P. Talpin. Implementation of the Typed Call-by-Value λ -Calculus using a Stack of Regions. In *Proceedings of the 21st Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 188–201, Portland, Oregon, Jan. 1994.
- [29] D. Walker and G. Morrisett. Alias Types for Recursive Data Structures. In *International Workshop on Types in Compilation*, Montreal, Canada, Sept. 2000.
- [30] D. Weise, 2001. Personal communication.
- [31] R. P. Wilson and M. S. Lam. Efficient Context-Sensitive Pointer Analysis for C Programs. In *Proceedings of the 1995 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 1–12, La Jolla, California, June 1995.
- [32] A. K. Wright. Typing References by Effect Inference. In B. Krieg-Brückner, editor, *4th European Symposium on Programming*, volume 582 of *Lecture Notes in Computer Science*, pages 473–491, Rennes, France, Feb. 1992. Springer-Verlag.
- [33] Z. Xu, T. Reps, and B. P. Miller. Typestate Checking of Machine Code. In D. Sands, editor, *10th European Symposium on Programming*, volume 2028 of *Lecture Notes in Computer Science*, pages 335–351, Genova, Italy, 2001. Springer-Verlag.