

go oracle: design

Alan Donovan
adonovan@google.com
August 25, 2013

The **go oracle** is a prototype source analysis tool, typically invoked by an editor, that answers questions about Go programs. This document describes its design.

There is also a [user manual](#).

[Motivation](#)
[Queries](#)
[Pointer analysis](#)
[Processing a query](#)
[Command-line interface](#)
 [Position](#)
 [Query mode](#)
 [Analysis scope](#)
[Example](#)
[Analysis libraries](#)
[Miscellaneous](#)
 [Well-formed inputs are required](#)
 [Complete Go sources are required for pointer analysis](#)
 [Status](#)
[Future work](#)
 [Features](#)
 [Performance optimisations](#)

Motivation

Programmers spend a great deal of time reading programs. Indeed, most of the time it takes to “write” a program is spent reading it, or more precisely, reading it and making logical deductions about what it does. Logic is, by definition, mechanical, and the goal of this work is to automate some of the deductions that Go programmers do, day in, day out in the course of their work so that, just as they no longer worry about whitespace and indentation (thanks to `gofmt`), they should no longer have to worry about a number of other mundane code comprehension tasks that a machine is better suited to carry out: locating definitions, ascertaining types of expressions, deducing the “implements” relation, computing method sets, finding callers/callees, jumping through channels, understanding aliasing.

Traditional IDEs such as Eclipse and IntelliJ support many of these kinds of queries (though not alias analysis---see below) but they are “omnivores”: heavyweight programs imposing a complex

theory of project organization on their users, and incorporating editors, build systems, debuggers, source control tools, code browsers and other utilities. We do not wish to pursue this approach for Go because the costs of supporting any specific IDE are very high and yet only the relatively few developers who use that IDE seriously are likely to benefit at all.

Instead, our approach with the oracle is to build a lightweight query system, not coupled to any particular environment, that can be called from almost any editor via a small amount of plumbing. The editor provides only the cursor position or selection and the query (e.g. “what’s the method-set of the selected expression?” or “who sends values on this channel?”) and the result, consisting of both source locations and informational messages, is printed in a compiler diagnostic format that the editor can mark up with hyperlinks in its usual way.

Queries

Here is a partial list of queries we plan to support.

- What is the type of this expression? What are its methods?
- What’s the value of this constant expression?
- Where is the definition of this identifier?
- What are the exported members of this imported package?
- What are the free variables of the selected block of code?
- What interfaces does this type satisfy?
- Which concrete types implement this interface?

And:

- What are the possible concrete types of this interface value?
- What are the possible callees of this dynamic call?
- What are the possible callers of this function?
- What objects might this pointer point to?
- Where are the corresponding sends/receives of this channel receive/send?
- Which statements could update this field/local/global/map/array/etc?
- Which functions might be called indirectly from this one?

In general, queries in the second group cannot be answered exactly because the answer may depend upon dynamic properties of the program; some static approximation is necessary. One approach is to use a **type-based approximation**. For example, we can approximate the set of potential callees at a dynamic function call by assuming, conservatively, that any function of the appropriate type may be called. For functions with unusual types, this gives good results, but for functions with common types this results in many spurious call edges; a tool that answered the query “what object might this *int pointer point to?” with the response “any int variable whose address is ever taken” would be unlikely to find many users.

Pointer analysis

A more precise approximation can be obtained using **pointer analysis**. This is a static analysis

technique in which all the statements of the entire program are analysed for their effect on *aliasing*, the relationship between pointers and the things they point to. There is a huge range of pointer analysis techniques in the literature, but generally, the more precise the technique (i.e. the fewer spurious results it delivers), the more expensive it is to compute.

Fortunately Go presents an attractive target for pointer analysis for several technical reasons:

- Go programs are almost completely type-safe.
Though some do use unsafe.Pointer conversions, these typically appear far less frequently than the analogous cast operations in C or C++.
- Go type hierarchies are shallow.
Compared to Java, in which very deep concrete class hierarchies are normal, implementation inheritance is impossible in Go, so all concrete types appear only on the lowest level of the type hierarchy.
- Go programs are not dynamically self-extending.
Go has no dynamic code loading/generating mechanism like `dlopen` or `mprotect (EXEC)` in C or `Class.forName` in Java, so it is feasible to analyse the complete source for the whole program statically.

Go programs are also smaller than C++ and Java programs; this may just be a consequence of the relative maturity of the language, but perhaps there is some effect of Go culture here too.

The oracle uses an **inclusion-based** pointer analysis, meaning that when it encounters a statement $y = x$, the analysis deduces that the set of things to which y may point *includes* (\subseteq) the set of things to which x may point. An alternative family of analyses is called **unification-based** because they pessimistically conclude that x and y may point to exactly the same set of things, i.e. they are *unified* (\equiv). Unification-based analyses scale well (linearly in the size of the program) but at the cost of poor precision. In contrast, solving an inclusion-based pointer analysis problem, which requires the computation of the transitive closure of a directed graph of inclusion constraints, requires cubic time and quadratic space in the size of the graph, so for many years it was considered impractical for large programs. However, recent advances in *presolver optimisation* have made it effectively linear by exploiting symmetry and redundancy in the graph; the solver is still $O(n^3)$, but the value of n that it sees is much smaller.

The pointer analysis is mostly **context-insensitive**, meaning that most functions are analyzed exactly once, so the analytic effect of each call to a function combines the effects of all calls to that function; only a few calls (e.g. built-ins, reflection, and some smaller functions) are treated context sensitively. Of course, this reduces precision, but generalized context-sensitive treatment is computationally very expensive.

As a prerequisite to pointer analysis, the program must first be converted from typed syntax trees into a simpler, more explicit **intermediate representation** (IR), as used by a compiler. We use a high-level **static single-assignment** (SSA) form IR in which the elements of all Go programs can be expressed using only about 30 basic instructions.

The pointer analysis is **flow-insensitive**, meaning that the order of the statements (sequencing, loops, conditionals) in the program has no effect. Such analyses are capable of answering *may alias* queries (“may P point to X?”) but not *must alias* queries (“must P point to X?”). While true flow-sensitive analysis is very expensive, the use of SSA-form IR gives a degree of flow sensitivity for free. For example, in the sequence

```
p = &x; *p = A  
p = &y; *p = B
```

SSA renaming would break p apart into two distinct variables, the first pointing only to x and the second only to y.

Processing a query

The steps required to process a query depend upon the mode of the query and also the selected user input. Typical queries involve the following steps:

- **Converting the filename and byte offset(s)** to an AST node.
This is done by walking the abstract syntax tree to find the smallest subtree that completely encloses the selection. Whitespace adjoining a subtree is treated as part of the subtree to make the tool more tolerant of sloppy selections.
- **Ascertaining the most appropriate AST node** for the query.
Certain queries require a particular type of node, such as a function call or a channel send or receive. The tool walks up (or sometimes down) the AST, starting at the selected node, to determine the appropriate node.

Many queries need only typed ASTs: those that connect definitions with references, enumerate the members of a package or method-set, or show the value of a constant expression. Such queries can be completely processed at this point. For queries needing pointer analysis, the following additional steps are required:

- **Locating the SSA value** corresponding to the source expression.
The SSA builder does some bookkeeping to record the relationship between syntax trees and SSA values (globals, parameters and instructions), analogous to a compiler maintaining debug information.
There may be zero, one or many SSA values for a given expression: zero if the expression is trivially dead code, one in the common case, or many if the function is analyzed multiple times due to context-sensitivity.
Each occurrence of a variable identifier is treated distinctly to obtain the local flow-sensitivity mentioned above.
- **Formulating and solving a pointer analysis problem.**
The pointer analysis query includes the set of SSA values of interest: for a “describe” query this is just the value of the selected expression; for a channel send/receive query, it would be all values ch for which there is a $<-ch$ or $ch<-x$ instruction anywhere in the program.
- **Displaying the result.**

Some bookkeeping is required to transform the results of pointer analysis into the desired output, such as a call graph, or a set of “may alias” facts about channel send and receive statements.

Command-line interface

An oracle command invocation consists of the following parts:

```
oracle -format=<plain|json> -pos=<filename>:<start>,<end> -mode=<mode> \
    <mainpkg> ...
```

Format

The `-format` flag specifies the desired output format: `json` or `plain` (the default). `plain` is a human-readable compiler diagnostic format; `json` is for consumption by other programs. The JSON interface is specified at go.tools/oracle/json/json.go; like everything else it is experimental and is likely to change.

Position

The `-pos` flag specifies a range of byte indexes within a particular source file, perhaps containing an identifier, an expression, or some other syntax of interest. This interval is known as the **selection**. This argument is trivial for any editor to construct, though admittedly not a natural or meaningful syntax for people.

Query mode

The `-mode` flag specifies the type of query to perform, one of the following:

- **callees**: show the possible call targets of the selected function call site.
- **callers**: show the possible callers of the function containing the selection.
- **callgraph**: show the entire callgraph of the program. (The selection is ignored.)
- **callstack**: show an arbitrary path from the root of the callgraph to the function containing the selection.
- **describe**: show various properties of the selected syntax: its syntactic kind, type, method set, constant value, point of definition. points-to set, etc, as appropriate.
- **freevars**: show the free variables of the selection.
- **implements**: show the ‘implements’ relation for all interfaces and concrete types defined in this package
- **peers**: show the set of possible sends/receives on the selected channel.
(The selection is a `<-` token.)

The grouping of queries is somewhat arbitrary. For example, the **peers** query, and others, could be merged into **describe** so that it is implicitly performed whenever the user ‘describes’ an expression of channel type.

Analysis scope

The `mainpkg` arguments specify one or more Go packages each containing a `main()` function. These packages, plus all packages transitively imported by them, form the scope of the pointer

analysis, typically including the applications being developed and their tests. We anticipate that users will configure the scope argument via their editor and that it will remain unchanged during an editing session.

Logically, the scope is a set of packages each defining a main function. Each such package can be:

- a command, i.e. a package defining an explicit main() function
e.g. code.google.com/p/go.tools/cmd/oracle
- a library with tests, i.e. a package defining a function named Test* called from the main package synthesized by 'go test'
e.g. fmt
- an ad-hoc main package comprised of a set of source files specified explicitly.
e.g. src/pkg/net/http/triv.go

The command-line syntax permits all three kinds of package to be specified but, due to current limitations of the type checker, at most one library-with-tests may be specified.

Example

Here is a concrete invocation of the oracle on a program in the Go standard library.

It corresponds to a query to **describe** the source selection shown in yellow:

```
triv.go:53:           io.Copy(buf, req.Body)
```

The output is shown in both formats, plain and json.

```
% oracle -format=... -pos=src/pkg/net/http/triv.go:#1042,#1050 -mode=describe \
src/pkg/net/http/triv.go
```

-format=plain

```
triv.go:53.20-53.23:      reference to var Body io.ReadCloser
request.go:110:2:          defined here
triv.go:53.16-53.23:      interface may contain these concrete types:
transfer.go:515:6:        *http.body, may point to:
transfer.go:333:18:        complit
transfer.go:338:17:        complit
transfer.go:343:18:        complit
-:                      *struct{*strings.Reader; io.Closer}, may point to:
server.go:1822:2:          complit
server.go:503:6:          *http.expectContinueReader, may point to:
server.go:1120:37:          complit
```

-format=json

```
{
  "mode": "describe",
  "describe": {
```

```

"desc": "identifier",
"pos": "/home/adonovan/go3/src/pkg/net/http/triv.go:53:20",
"detail": "value",
"value": {
    "type": "io.ReadCloser",
    "objpos": "/home/adonovan/go/src/pkg/net/http/request.go:110:2",
    "pts": [
        {
            "type": "*http.body",
            "namepos": "/home/adonovan/go/src/pkg/net/http/transfer.go:515:6",
            "labels": [
                {
                    "pos": "/home/adonovan/go/src/pkg/net/http/transfer.go:333:18",
                    "desc": "complit"
                },
                {
                    "pos": "/home/adonovan/go/src/pkg/net/http/transfer.go:338:17",
                    "desc": "complit"
                },
                {
                    "pos": "/home/adonovan/go/src/pkg/net/http/transfer.go:343:18",
                    "desc": "complit"
                }
            ]
        },
        {
            "type": "*http.expectContinueReader",
            "namepos": "/home/adonovan/go/src/pkg/net/http/server.go:503:6",
            "labels": [
                {
                    "pos": "/home/adonovan/go/src/pkg/net/http/server.go:1120:37",
                    "desc": "complit"
                }
            ]
        },
        {
            "type": "*struct{*strings.Reader; io.Closer}",
            "labels": [
                {
                    "pos": "/home/adonovan/go/src/pkg/net/http/server.go:1822:2",
                    "desc": "complit"
                }
            ]
        }
    ],
}

```

In response to this query, the oracle has provided the type of the expression, the location of the definition of the struct field, the list of concrete types that it (an interface) may contain, and for each of those concrete types, all of which are pointers, the set of objects to which it may point, with source locations where available. (`complit` indicates the object allocated by a composite

literal.)

In the plain format output, only the base names of the the resulting filenames are shown here, for brevity.

In both formats, the tool also prints a large number of pointer analysis warnings, not shown. Most arise from intrinsic functions that the analysis does not yet know about, and will disappear as support for the intrinsics is added. We may need to present them in a less intrusive manner in the user interface though.

The value of the -pos flag is easily is constructed by an editor but not by a human user. A minority of queries do not require a position (e.g. **callgraph**) or could easily be altered to take their input in a different form (e.g. **describe** applied to a package name), but for most queries there is simply no convenient textual way to indicate a specific syntax tree.

The current command-line interface is undoubtedly crude but should suffice to display the potential of static analysis to assist programmers with code comprehension. No aspect of the user interface should be considered set in stone, and we welcome feedback from early users on how it could be improved.

Analysis libraries

The oracle uses the following libraries in the [golang.org repo](#):

- | | |
|---|-----------------------|
| • go/{token,scanner,ast,parser} | parser |
| • code.google.com/p/go.tools/go/types | type checker |
| • code.google.com/p/go.tools/importer | source package loader |
| • code.google.com/p/go.tools/ssa | SSA IR |
| • code.google.com/p/go.tools/pointer | pointer analysis |
| • code.google.com/p/go.tools/oracle | oracle library |
| • code.google.com/p/go.tools/cmd/oracle | oracle command |

The SSA and pointer analysis libraries were developed for and with the oracle, but we plan to build or extend other applications to make use of them, for example

- improved documentation tools (c.f. godoc)
- improved bug-finding tools (c.f. go vet)
- refactoring tools

Restrictions

Well-formed inputs are required

Since the oracle is based on the Go standard `go/parser` package, which rejects malformed

inputs just like a compiler does, and since the SSA builder requires a well-formed Go program as its input, the oracle is not forgiving of ill-formed programs. Tools such as Eclipse use fault-tolerant parsers that repair damage as best they can so they can always make progress. We may experiment with such an approach, but in the meantime, users wanting best results from the oracle will need to adapt their coding workflow to ensure that their programs are compilable at all times. An unfortunate consequence of this limitation is that name completion while typing new code cannot currently be provided by the oracle.

Complete Go sources are required for pointer analysis

Pointer analysis requires that the complete Go source of the program is available. The effects of functions written in assembly or C cannot be observed by the analysis, causing imprecise or unsound results, although often these problems only appear in the vicinity of the native code. We plan to provide a mechanism for users to summarize (in Go source) the aliasing effects of their native code to improve precision where needed.

go/build package layout conventions are assumed

The oracle uses the `go/build` package to find the source files for a given package. However, the ‘`go test`’ command has significant additional logic for constructing test packages and the oracle cannot faithfully reproduce this in all cases (yet). Furthermore, some proprietary build systems (such as Google’s) have a slightly different package organization requiring a different algorithm to locate the source files given an import path. The oracle cannot be run in such environments (yet).

Status

As of September 1, 2013, the oracle is available to adventurous users for evaluation; we welcome feedback and bug reports.

Known bugs

- All the known bugs of the `go/types` package.
In particular: cyclic package dependencies are not yet detected and reported and will cause the program to get stuck.
- All the known bugs of the `go.tools/pointer` package
 - Neither `unsafe.Pointer` nor reflection are treated soundly by the pointer analysis.
 - See also: [go.tools/pointer/TODO](#)
- The algorithm that converts source positions to AST nodes yields some surprising results.
- The Emacs markup is too slow, especially in **callgraph** mode
- **freevars** prints some package-level names
- **peers**
 - permit querying from a `makechan`, `close()`, `for...range`, or reflective op
 - also report aliasing `reflect.{Send,Recv,Close}` and `close()` operations

Future work

Features

- Alternative UI: rather than `-format=plain|json|...` etc, consider having the user specify a template that runs over (say) the JSON structures to generate the output in the desired form. See ‘go list’ for inspiration.
- Include complete [start...end) extent information in JSON output, not just a position, where available. (One subtlety is that sometimes a position is better than an extent. Consider the expression $e_1 + e_2$ where e_1 and e_2 are complex expressions: if you want to indicate the addition operator, the position of the “+” token is much more helpful than the entire extent of the expression.)
- Allow the analysis scope to include multiple test packages at once.
(Requires go/types changes.)
- Allow the tool to run ‘describe’ queries providing type information (but no pointer analysis) even on packages that are outside the analysis scope.
- Add support for refactorings such as extract/inline method, renamings, add/remove/reorder parameters, etc.
- Add support for other editors: Vim, Acme, Sublime, Eclipse, etc. (Volunteers wanted.)
- Improve the user interface to permit single-button traversal around the callgraph.
- Provide a mechanism to allow users to summaries the aliasing effects of their native code, for improved analysis precision.
- Add an **updates** query: “which statements can update this (addressable) expression?”

Performance optimisations

- Ideally the tool should respond almost instantaneously to most queries, and we believe this is quite feasible for all but the largest programs, without resorting to stateful or long-lived analysis processes.
- Implement pre-solver optimisations in the pointer analysis; ideally its running time should be barely noticeable.
- Parallelize the type-checking of function bodies.
- Cache the results of callgraph-related queries across invocations in the editor so that navigation is instantaneous.