

# Abstraction in Technical Computing

by

Jeffrey Werner Bezanson

Submitted to the Department of Electrical Engineering and Computer  
Science

in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

February 2015

© Massachusetts Institute of Technology 2015. All rights reserved.

Author .....  
Department of Electrical Engineering and Computer Science  
January 9, 2015

Certified by .....  
Alan Edelman  
Professor  
Thesis Supervisor

Accepted by .....  
Leslie Kolodziejski  
Chairman, Department Committee on Graduate Students

# Abstraction in Technical Computing

by

Jeffrey Werner Bezanson

Submitted to the Department of Electrical Engineering and Computer Science  
on January 9, 2015, in partial fulfillment of the  
requirements for the degree of  
Doctor of Philosophy

## Abstract

Array-based programming environments are popular for scientific and technical computing. These systems consist of built-in function libraries paired with high-level languages for interaction. Although the libraries perform well, it is widely believed that scripting in these languages is necessarily slow, and that only heroic feats of engineering can at best partially ameliorate this problem.

In this thesis I argue that what is really needed is a more coherent structure for this functionality. To find one, we must ask what technical computing is really about. I suggest that this kind of programming is characterized by an emphasis on operator complexity and code specialization, and that a language can be designed to better fit these requirements.

The key idea is to integrate code *selection* with code *specialization*, using generic functions and data-flow type inference. Systems like these can suffer from inefficient compilation, or from uncertainty about what exactly to specialize on. I show that dispatch on structured type tags helps address these problems. The resulting language, Julia, achieves a Quine-style “explication by elimination” of many of the productive features technical computing users expect.

Thesis Supervisor: Alan Edelman  
Title: Professor

## Acknowledgments

thank Tim Holy for frequently and spontaneously writing explanations that I found myself wanting to copy and paste into this writeup.

# Contents

<b>1</b>	<b>Introduction</b>	<b>7</b>
1.1	The technical computing problem . . . . .	8
1.2	Proposed solution . . . . .	8
1.3	Contributions . . . . .	10
<b>2</b>	<b>Problems in technical computing</b>	<b>11</b>
2.1	What is a technical computing environment? . . . . .	11
2.2	The software stack is too complex . . . . .	11
2.3	Why dynamic typing? . . . . .	12
2.3.1	Operational reasoning . . . . .	12
2.3.2	I/O . . . . .	13
2.3.3	Flat metadata hierarchy . . . . .	13
2.3.4	Mismatches with mathematical abstractions . . . . .	13
2.4	Bad tradeoffs in current designs . . . . .	14
2.4.1	What needs to be built in? . . . . .	15
2.5	Social phenomena . . . . .	16
<b>3</b>	<b>Available technology</b>	<b>18</b>
3.1	The language design space . . . . .	18
3.1.1	Classes vs. functions . . . . .	19
3.1.2	Separate compilation vs. performance . . . . .	19
3.1.3	Parametric vs. ad-hoc polymorphism . . . . .	19
3.1.4	Static checking vs. flexibility . . . . .	19

3.1.5	Modularity vs. large vocabularies . . . . .	19
3.1.6	Eliminating vs. embracing tagged data . . . . .	19
3.1.7	Data hiding vs. explicit memory layout . . . . .	19
3.1.8	Dynamic dispatch is key . . . . .	19
3.2	Multiple dispatch . . . . .	25
3.2.1	Symbolic programming . . . . .	25
3.3	Domain theory . . . . .	26
3.4	Dynamic type inference . . . . .	27
3.4.1	Objections to dynamic typing . . . . .	27
<b>4</b>	<b>The Julia approach</b>	<b>29</b>
4.1	Language description . . . . .	29
4.1.1	Type system . . . . .	29
4.1.2	Data model . . . . .	29
4.1.3	Dispatch mechanism . . . . .	29
4.2	Performance model . . . . .	29
4.2.1	Type inference . . . . .	29
4.2.2	Specialization . . . . .	29
4.3	Explication through elimination . . . . .	29
4.3.1	Conversion and other basic operations . . . . .	29
4.3.2	Numeric types and embeddings . . . . .	30
4.3.3	Multidimensional array indexing . . . . .	33
4.3.4	Array views . . . . .	37
4.3.5	Units . . . . .	37
4.3.6	Even more elimination? . . . . .	37
<b>5</b>	<b>Case studies</b>	<b>38</b>
5.1	Numerical linear algebra . . . . .	38
5.2	Grid.jl, computational libraries . . . . .	38
5.3	JUMP . . . . .	38
5.4	ACAS . . . . .	38

5.5	Parallel computing, e.g. parallel prefix . . . . .	38
5.6	IJulia . . . . .	38
<b>6</b>	<b>Conclusion</b>	<b>39</b>
	<b>Bibliography</b>	<b>41</b>

# Chapter 1

## Introduction

Much of the history of programming languages has been about increasing abstraction and generalization. Today we take it for granted that a single suitably powerful language could be used for nearly any programming task, were it not for pragmatic factors beyond our control. Over time, many special-purpose languages have been subsumed by these more powerful and general languages. The field of technical computing — programming for applied math and the sciences — is something of a holdout from this trend, with its own purpose-built languages still dominating. There are probably many reasons for this, including historical, social and technical.

In this work we will focus on the technical reasons for this, and the problems that the current state of affairs causes.

What do we mean by a *language* for scientific computing? The prevailing answer is that it is about numeric arrays, perhaps matrices specifically, and numerical libraries that use them. Instead we intend to argue that, at a deeper level, it is about certain kinds of flexibility, particularly flexibility in the behavior of key operators and functions. This flexibility is needed to maximize composability and reusability in scientific code. Without it, certain programs may be easy to write today, but changing functional and performance requirements can become difficult to meet.

It is clear that scientific and technical users have their own particular needs. So far there have been two paths available for fulfilling these needs: using a mainstream general-purpose language like C++ or Java, or designing a new language like R or

Mainstream PL	Technical computing
classes, single dispatch	complex operators
separate compilation	performance, inlining
parametric polymorphism	ad-hoc polymorphism, extensibility
static checking	experimental computing
modularity, encapsulation	large vocabularies
eliminating tags	self-describing data, acceptance of tags
data hiding	explicit memory layout

Figure 1-1: Priorities of mainstream object-oriented and functional programming language research and implementation compared to those of the technical computing domain.

MATLAB. Designing your own language seems like a drastic thing to do, and yet it is unusually common in technical computing (see gang of 40).

## 1.1 The technical computing problem

Figure 1.1 compares the general design priorities of mainstream programming languages to those of technical computing languages. The priorities in each row are not necessarily opposites or even mutually exclusive, but rather are a matter of emphasis. For example, it is certainly possible to have both parametric and ad-hoc polymorphism within the same language, but syntax, recommended idioms, and the design of the standard library will tend to emphasize one or the other.

It is striking how different these priorities are. We believe these technical factors have contributed significantly to the persistence of specialized environments in this area.

## 1.2 Proposed solution

Solution: integrate code specialization and code selection

Why integrate code specialization and code selection?

- specialization requires selection anyway



- simpler system
- can sometimes collapse 2 layers of dispatch into 1
- can replace library code with a code generator without either changing client code *\*or\** any extra overhead

Extensive code specialization is a key feature of technical computing. “what to specialize on” has been an open problem. our types are a possible solution to this for two reasons:

1. you can tune the amount of information they contain
2. everybody agrees to use them, which helps ensure that specializing on the types will actually do something useful.

The key ingredients:

1. Self-describing data model aware of memory layout
2. Type tags with nested structure
3. A fully-connected type tree
4. Dynamic multiple dispatch over all types
5. Dataflow type inference
6. Automatic code specialization

This list of features may appear somewhat ad-hoc. However, they turn out to be remarkably strongly coupled, and deeply constrained by our ultimate goal. Each of these features has appeared in some form before, but never in a way that fully solves the problems described here.

Challenges of this approach (why has this not been done before?)

## 1.3 Contributions

“One of the most fruitful techniques of language analysis is explication through elimination. The basic idea is that one explains a linguistic feature by showing how one could do without it.” [1]

A main contribution of this thesis is the application of this approach to features of technical computing environments that have not been subject to such analysis before.

# Chapter 2

## Problems in technical computing

At a high level, the experiences and thoughts that give rise to our design decisions.

How should a computer scientist approach this space? We might try to maximize performance. Or it is also easy to make unfounded assumptions about “what users want”. But instead we should study the real world, see what is happening and figure out how to steer it in a better direction.

Hypothesis: people don’t know what they want. It’s also hard to predict what people will want in the future. We need, in the words of Gerry Sussman, systems adaptable to uses not imagined by their designers.

### 2.1 What is a technical computing environment?

This question has not really been answered.

Views of this are strongly shaped by what systems happen to exist, and what people were exposed to as they learned to program.

Technical computing software has been designed haphazardly. Each system is a pile of features taken without argument.

### 2.2 The software stack is too complex

Collapsing abstraction layers

- for numerical debuggability in particular - debuggability in general - no time spent worrying about binding time - you will build a dynamic dispatch layer anyway, so build it in

How much of the past 30 years of handwritten Matlab internals can be autogenerated with a compiler? (A lot)

Can now get past traditional classifications of language boundaries - interpreted vs. compiled, high-level vs. low-level, procedural vs. imperative vs. functional, etc.

language performance psychology: if your language doesn't directly support efficient machine data types, users will rewrite their code in C in order to get them, and then be happy with the result (though not with the process).

so julia is an all-levels language

## 2.3 Why dynamic typing?

Mathematical abstractions often frustrate our attempts to represent them on a computer. For example, mathematicians can move instinctively between isomorphic objects such as scalars and 1x1 matrices, but most programming languages would prefer to represent scalars and matrices quite differently. A system might be able to convert automatically from one to the other, but it cannot always know which one we \*meant\*.

### 2.3.1 Operational reasoning

people tend to think about programs operationally, i.e. what it \*does\* when it runs. for example writing `if false code end` the code does not "occur" and therefore does not need to be valid

there is less to learn. with static languages you have to learn what happens at both compile time and run time, when only run time really matters.

there is a desire to parameterize as much as possible. functions accept parameters, so function calling ought to be sufficient to express any desired parameterization.

### 2.3.2 I/O

inevitably there is a need to refer to a datatype at run time. the best example is file I/O, where you might want to say "read 500 double-precision numbers from file X". in static languages the syntax and identifiers used to specify such run-time types must be different from those used to specify static types. in C you see defined constants such as `DATATYPE_DOUBLE`.

### 2.3.3 Flat metadata hierarchy

static types are approximations of dynamic types, so languages with static types inevitably assign two types to a location (both a static type and a dynamic type) where one would do. in some languages, like C++, the desire for performance or ease of implementation leads the compiler to make some decisions based on static types. this is confusing. if type declarations can be omitted, as in a type-inferred language, the situation is even worse since the static type of a value might not be apparent.

### 2.3.4 Mismatches with mathematical abstractions

programs, in general, deal with values of widely varying disjoint types: functions, numbers, lists, network sockets, etc. type systems are good at sorting out values of these different types. however, in mathematical code most values are numbers. numerical properties (such as positive, negative, even, odd, greater than 1, etc.) are what matter, and these are highly dynamic. the lattices involved are generally not of finite height.

different number representations exist only for efficiency, and are often incidental to the meaning of a piece of code. for example people want " $y = \sqrt{x}$ " to compute the square root of  $x$  whether it is positive or negative, and give a real or complex result accordingly. as another example, in many cases it is convenient not to need to worry about integer overflow.

multiple common features underlie mathematical objects of different types (e.g. numbers, sets, matrices). in some cases it makes sense to consider numbers and

matrices as the same kind of thing, and in other cases it doesn't matter. A given type system is likely not to have anticipated the particular common features that matter to your program, making it more difficult to express an idea. A concrete example is the matlab fragment `if condition idx = ':' else idx = 1 end` where we want to select either an entire dimension or the first position alone. The `:` and `1` are both indexes in this context, though they would be of disjoint types in most programming languages.

## 2.4 Bad tradeoffs in current designs

- Ducking the issue of typing altogether - why this isn't possible IRL

Case in point: Python and NumPy.

NumPy tried to work within the confines of Python, but has more or less failed for technical and political reasons. (Technical: hard to work with types in language that deliberately tries to obscure them from the user. Political: unwillingness to fix broken package distribution system.) Consequences: Numba and Numba.lang, Anaconda. Essentially reinventing types in "Python".

This phenomenon is increasingly noticed in other domains, particularly JavaScript and web programming. Modern JavaScript implementations are quite fast, but Google's Dart language is based on the premise that we could have a web language that is even faster, and offers more productivity as well. How so? Because Dart's designers observed that JavaScript programmers in practice often write code that could be defined using traditional OO classes, but the language does not support them.

dictionaries for everything (python, js) is the wrong default. almost every type somebody wants has a fixed number of fields with fixed types.

Python is often described as a good glue language. This means it is effectively used as an interface standard, a kind of extended C ABI that makes it easier for libraries to interoperate, and easier for users to access those libraries. Something as straightforward as providing a standard N-dimensional numeric array class (NumPy), which does not exist at the level of C, goes a long way.

However, we wish to point out that in this picture, Python is not doing as much

work as it might first appear. Python does not make it easier to implement the functionality inside NumPy, or other “under the hood” scientific libraries. In many cases it creates more work, through the need to write wrappers and interfaces for native code.

Merely being “dynamic” (e.g. Python) should not be considered the gold standard of flexibility. Although these systems permit tricks that can solve otherwise difficult programming problems, this is not always the kind of flexibility that is needed. When faced with the need to describe many functions with elaborate behavior and many cases, one does not primarily need permissiveness, but rather powerful and descriptive organizing principles.

How does Julia break through the 4x-slower-than-C performance barrier?

bits types and parameters – can have abstract types whose size is known unconditionally

### 2.4.1 What needs to be built in?

On “built-in”, why built-in is bad. Built-in-ness often conflates two aspects:

1. A feature being readily available and agreed-on by all language users
2. A feature tightly coupled to the rest of the system

(2) implies (1), but not the other way around. (2) is the only technically interesting item, since the other can be addressed e.g. just by including a library in the standard software distribution. Many technical computing languages have done a large amount of (2) while justifying it with point (1).

While a large part of our motivation is to move more decisions and functionality into libraries, it is equally important to identify what *\*must\** be part of a language for the system to be successful. We believe that large amounts of functionality can be provided by add-ons, but that certain key features absolutely cannot be. Past failures to properly classify features this way have caused a lot of undue pain.

First, performance cannot be an add-on. If some users have a fast version of a language and others have a slow version (with the difference being an order of magnitude

or more), library writers cannot be sure whether users will find their code fast enough to be useful. How are we to teach people to program in the language? Courses on MATLAB programming emphasize vectorized code, but if not all implementations had this performance characteristic the curriculum would become confused.

Psychologically, it may be difficult to accept a “non standard” extension that changes a language so fundamentally. There is a nagging, though perhaps totally unfounded, perception that something subtle may break. If indeed a bug arises due to the use of such an extension, a user is likely to conclude that the extension is dangerous or broken and stop using it. If, on the other hand, a bug arises due to a language’s standard optimizing compiler, the user will simply file a bug report, then find a way to work around the problem.

Adding a JIT compiler to a language also requires acceptance of detriments like compilation pauses and pages with RWX permissions. In some cases this may lead to use of the extension being disallowed, perhaps for security reasons.

Type systems similarly fail when provided as optional extensions. Library writers face the same kinds of problems as with performance add-ons. Should I use type annotations in my library?

Dynamic dispatch mechanisms also make especially poor add-ons. Of course, every program makes decisions at run time, and so implements its own “dispatch” to some extent. But these behaviors are inextensible; if language users do not agree on a reusable dispatch framework their code will not be composable.

## 2.5 Social phenomena

Programming languages are observed to have strong network effects, and the difficulty of getting new languages adopted is well known. However based on [socio PLT paper] we believe this doesn’t have to be the case. The formula of improving or redesigning general-purposes languages to be more appealing to domain experts might solve the problem. That way the new system has immediate appeal for at least some users, without the worry that a different tool will be needed as soon as requirements change



slightly.

barrier to contributing

Software business is based on imposing restrictions

Debates about what abstractions mean – AbstractMatrix we thought we knew what it meant, but what about something like SymTridiagonal? It can implement most of what a dense matrix has, but it can't obey every invariant.

PackedQR

# Chapter 3

## Available technology

### 3.1 The language design space

It is helpful to begin with a rather coarse classification of programming languages, according to how expressive their type systems are, and whether their type systems are dynamic or static:

	More types	Fewer types
Dynamic	Dylan, Julia	Scheme, Python, MATLAB
Static	ML, Haskell, Scala	C

The lower right corner tend to contain older languages. The upper right corner contains many popular “dynamic” languages. The lower left corner contains many modern languages resulting from research on static type systems. We are most interested in the upper left corner, which is notable for being rather sparsely populated. It has been generally believed that dynamic languages do not “need” types, or that there is no point in talking about types if they are not going to be checked at compile time. These views have some merit, but as a result the top-left corner of this design space has been seriously under-explored.

### **3.1.1 Classes vs. functions**

### **3.1.2 Separate compilation vs. performance**

### **3.1.3 Parametric vs. ad-hoc polymorphism**

The idea of specialization unites parametric and ad-hoc polymorphism. Beginning with a parametrically-polymorphic function, one can imagine a compiler specializing it for various cases, i.e. certain concrete argument values. These specializations would be stored in a lookup table, for use either at compile time or at run time.

Next we can imagine that the compiler might not optimally specialize certain definitions, and that the programmer would like to provide hints or hand-written implementations to speed up special cases. For example, imagine a function that traverses a generic array. A compiler-generated specialization might inline a specific array types's indexing operations, but a human might further realize that the loop order should be switched for certain arrays types based on their storage order.

However, if we are going to let a human specialize a function for performance, we might as well allow them to specialize it for some other reason, including entirely different behavior. But at this point separate ad-hoc polymorphism is no longer necessary; we are using a single overloading feature for everything.

### **3.1.4 Static checking vs. flexibility**

### **3.1.5 Modularity vs. large vocabularies**

### **3.1.6 Eliminating vs. embracing tagged data**

### **3.1.7 Data hiding vs. explicit memory layout**

### **3.1.8 Dynamic dispatch is key**

It would be unpleasant if every piece of every program we wrote were forced to do only one specific task. Every time we wanted to do something slightly different, we'd have to write a different program. But if a language allows the same program element

to do different things at different times, we can write whole classes of programs at once. This kind of capability is one of the main reasons *object-oriented* programming is popular: it provides a way to automatically select different behaviors according to some structured criteria.

In class-based OO there is essentially *no way* to create an operation that dispatches on existing types (the expression problem). This clearly does not match technical computing, where most programs deal with the same few types (e.g. number, array), and might sensibly want to write new operations that dispatch on them.

We use the non-standard term “criteria” deliberately, in order to clarify our point of view, which is independent of any particular object system.

The sophistication of the available “selection criteria” account for a large part of the perceived “power” or leverage provided by a language. In fact it is possible to illustrate a hierarchy of such mechanisms. As an example, consider a simple simulation, and how it can be written under a series of increasingly powerful paradigms. First, written-out imperative code:

```
while running
  for a in animals
    b = nearby_animal(a)
    if a isa Rabbit
      if b isa Wolf then run(a)
      if b isa Rabbit then mate(a,b)
    else if a isa Wolf
      if b isa Rabbit then eat(a,b)
      if b isa Wolf then follow(a,b)
    end
  end
end
end
```

We can see how this would get tedious as we add more kinds of animals and more behaviors. Another problem is that the animal behavior is implemented directly

inside the control loop, so it is hard to see what parts are simulation control logic and what parts are animal behavior. Adding a simple object system leads to a nicer implementation <sup>1</sup>:

```
class Rabbit
  method encounter(b)
    if b isa Wolf then run()
    if b isa Rabbit then mate(b)
  end
end

class Wolf
  method encounter(b)
    if b isa Rabbit then eat(b)
    if b isa Wolf then follow(b)
  end
end

while running
  for a in animals
    b = nearby_animal(a)
    a.encounter(b)
  end
end
```

Here all of the simulation's animal behavior has been essentially compressed into a single program point: `a.encounter(b)` leads to all of the behavior by selecting an implementation based on the first argument, `a`. This kind of criterion is essentially indexed lookup; we can imagine that `a` is simply an integer index into a table of operations.

---

<sup>1</sup>A perennial problem with simple examples is that better implementations often make the code longer.

The next enhancement to “selection criteria” adds a hierarchy of behaviors, to provide further opportunities to avoid repetition:

```
abstract class Animal
  method nearby()
    # search within some radius
  end
end

class Rabbit <: Animal
  method encounter(b: Animal)
    if b isa Wolf then run()
    if b isa Rabbit then mate(b)
  end
end

class Wolf <: Animal
  method encounter(b: Animal)
    if b isa Rabbit then eat(b)
    if b isa Wolf then follow(b)
  end
end

while running
  for a in animals
    b = a.nearby()
    a.encounter(b)
  end
end
```

We are still essentially doing table lookup, but the tables have more structure:

every `Animal` has the `nearby` method, and can inherit a general-purpose implementation.

This brings us roughly to the level of most popular object-oriented languages. But in this example still more can be done. Notice that in the first step to objects we replaced one level of `if` statements with method lookup. However, inside of these methods a structured set of `if` statements remains. We can replace these by adding another level of dispatch.

```
class Rabbit <: Animal
  method encounter(b: Wolf) = run()
  method encounter(b: Rabbit) = mate(b)
end
```

```
class Wolf <: Animal
  method encounter(b: Rabbit) = eat(b)
  method encounter(b: Wolf) = follow(b)
end
```

We now have a *double dispatch* system, where a method call uses two lookups, first on the first argument and then on the second argument.

This syntax might be considered a bit nicer, but the design clearly begs a question: why is  $n = 2$  special? It isn't, and we could clearly consider even more method arguments as part of dispatch. But at that point, why is the first argument special? Why separate methods in a special way based on the first argument? It seems arbitrary, and indeed we can remove the special treatment:

```
abstract class Animal
end
```

```
class Rabbit <: Animal
end
```

```

class Wolf <: Animal
end

nearby(a: Animal) = # search
encounter(a: Rabbit, b: Wolf) = run(a)
encounter(a: Rabbit, b: Rabbit) = mate(a,b)
encounter(a: Wolf, b: Rabbit) = eat(a, b)
encounter(a: Wolf, b: Wolf) = follow(a, b)

while running
  for a in animals
    b = nearby(a)
    encounter(a, b)
  end
end
end

```

Here we made two major changes: the methods have been moved “outside” of any classes, and all arguments are listed explicitly. This change has fairly significant implications. Since methods no longer need to be “inside” classes, there is no syntactic limit to where definitions may appear. Now it is easier to add new methods after a class has been defined. Methods also now naturally operate on combinations of objects, not single objects.

The shift to thinking about combinations of objects is fairly revolutionary. Many interesting properties only apply to combinations of objects, and not individuals. We are also now free to think of more exotic kinds of combinations.

We can define a method on *any number* of objects:

```

encounter(ws: Wolf...) = pack(ws)

encounter{T<:Animal}(a: T, b: T) = mate(a, b)

```



## 3.2 Multiple dispatch

CHART of different multiple dispatches. classify them

Why “just overloading” is not enough. You intercept every operation and rewrite it, kind of an escape hatch for a language you don’t like. If somebody else also tries to do this, there won’t necessarily be any coherence.

A helpful way to classify languages with some kind of generic programming support is to look at which language constructs are generic. For example, in C++ the syntax `object.method(x)` is generic: a programmer can get it to do different things by supplying different values for `object`. The C++ syntax `f(a,b)` or `a+b` is usually not generic, but can be overloaded by supplying definitions for different argument types. However, this overloading only uses compile-time types (no run-time information), and so is essentially a form of renaming — it can be implemented by renaming each definition with a unique name, and replacing overloaded calls with an appropriate name based on compile-time types. This is a much weaker form of generic programming, since different behaviors cannot be obtained by passing different values at run time.

For this reason, C++ programs are not fully generic: it is difficult to substitute new behaviors for every part of a program. Some languages take generic programming much further. For example, in MATLAB, *every* function is effectively generic, and can be overloaded by new classes. This enables a programmer to “intercept” every function call (and therefore, essentially everything a program does)

### 3.2.1 Symbolic programming

Systems based on symbolic rewrite rules arguably occupy a further tier of dispatch sophistication. In these systems, you can dispatch on essentially anything, including arbitrary values and structures. These systems are typically powerful enough to concisely define the kinds of behaviors we are interested in.

However, symbolic programming lacks data abstraction: the concrete representations of values are exposed to the dispatch system (e.g. there is no difference between

being a list and being something represented as a list).

Patterns are very powerful, but the tradeoff is that there is not necessarily a useful relationship between what your program does and what a static analysis (based on a finite-height partial order over patterns) can discover. Maybe julia could be considered a sweet spot somewhere in between.

### 3.3 Domain theory

In the 1960s Dana Scott asked how to assign meanings to programs, which otherwise just appear to be lists of symbols. For example, given a program computing the factorial function, we want a process by which we can assign the meaning “factorial” to it.

This is about modeling the behavior of a program without running it.

The idea of analyzing computer programs began in earnest in the 1960s with Dana Scott’s invention of domain theory. A “domain” in this theory is a partial order of sets of values that a program might manipulate. Domain theory models computation as follows: a program starts with no information, the lowest point in the partial order (“bottom”). Computation steps accumulate information, gradually moving higher through the order. The advantage of this model, in essence, is that it provides a way to think about the meaning of a program without running the program. Even the “result” of a non-terminating program has a representation — the bottom element. Other elements of the partial order might refer to intermediate results.

Domain theory gave rise to the study of denotational semantics and the design of type systems. However, the original theory is quite general and invites us to invent any domains we wish for any sort of language. For example, given a program that outputs an integer, we might decide that we only care whether this integer is even or odd. Then our posets are the even and odd integers, and we will classify operations in the program according to whether they are evenness-preserving and so on. It should be clear that this sort of analysis, while clearly related to type systems, is fairly different from what most programmers think of as type checking.

## 3.4 Dynamic type inference

This is exciting because the type system can do useful computational work for you, instead of only checking the work done by the rest of the program.

Joins of quantified types

“Diagonal dispatch”

(T,T) (Int, Real)

$\backslash(\text{Int}, \text{Int})/$

Avoid hard coding logic into the compiler.

Exposing to the user what is part of the language anyway and not just hiding it in the compiler.

Data flow analysis has to actually work

Taking things out of the language that break dataflow analysis

Local reasoning

You are imposing a type system anyway, even implicitly perhaps, so might as well make that from the get-go.

### 3.4.1 Objections to dynamic typing

it may be that the “power” of a language is equal to the complexity of the criteria used by the language’s run-time dispatch mechanisms.

pre-OO: pointer indirection OO: single dispatch, class hierarchies

Haskell: without typeclasses, a beautiful language, but one that nobody would use.

I claim that compile-time abstractions do not count. The problem is that at some point you have to \*actually run the program\*. Specifying the behavior of a program is hard; this is where we need help from the language.

Scripting languages are often defended using the observation that most code is not performance-critical. However, this fact does not mean that it’s ok for a language design to ignore performance, nor does it mean that performance should be the default priority of every line of code. Rather it means that the default should be convenience

and safety, with a path to performance easily in reach for when it's needed.

# Chapter 4

## The Julia approach

### 4.1 Language description

#### 4.1.1 Type system

#### 4.1.2 Data model

#### 4.1.3 Dispatch mechanism

### 4.2 Performance model

#### 4.2.1 Type inference

#### 4.2.2 Specialization

### 4.3 Explication through elimination

#### 4.3.1 Conversion and other basic operations

This section will illustrate how we implement key features of technical computing systems using our methodology.

- The abstractions of equality and comparison. Different equivalence classes between

is/==, isequal and ==  
- Numeric vs lexicographic ordering?  
cmp, lexcmp, vs isless, j

### 4.3.2 Numeric types and embeddings

We might prefer “number” to be a single, concrete concept, but the history of mathematics has seen the concept extended many times, from integers to rationals to reals, and then to complex, quaternion, and more. These constructions tend to follow a pattern: a new set of numbers is constructed around a subset isomorphic to an existing set of numbers. For example, the reals are isomorphic to the complex numbers with zero imaginary part.

Human beings happen to be good at equating and moving between isomorphic sets, so it is easy to imagine that the reals and complexes with zero imaginary part are one and the same. But a computer forces us to be specific, and admit that a real number is not complex, and a complex number is not real. And yet the close relationship between them is too compelling not to model in a computer somehow. Here we have a numerical analog to the famous “circle and ellipse” problem in object-oriented programming: the set of circles is isomorphic to the set of ellipses with equal axes, yet neither “is a” relationship in a class hierarchy seems fully correct. An ellipse is not a circle, and in general a circle cannot serve as an ellipse (for example, the set of circles is not closed under the same operations that the set of ellipses is, so a program written for ellipses might not work on circles). This problem implies that a single built-in type hierarchy is not sufficient: we want to model custom \*kinds\* of relationships between types (e.g. “can be embedded in” in addition to “is a”).

Two further problems should also be kept in mind. First, the natural isomorphisms between sets of numbers might not be isomorphisms on a real computer. For example, due to the behavior of floating-point arithmetic, an operation on complex numbers with zero imaginary part might not give an answer equal to the same operation on real numbers. Second, the contexts that demand use of one type of number or another are often not easily described by type systems. The classic example is square root (`sqrt`),

whose result is complex for negative arguments. Including a number’s sign in its type is a possibility, but this quickly gets out of hand — should a type system attempt to prove a matrix symmetric before we compute its eigenvalues? While we cannot offer a once-and-for-all solution to these problems, we will show how the flexibility of our proposed mechanism is useful for addressing them.

## Applications

*Ranges* illustrate an interesting application of type promotion. A range data type, notated `a:s:b`, represents a sequence of values starting at `a` and ending at `b`, with a distance of `s` between elements (internally, this notation is translated to `colon(a, s, b)`). Ranges seem simple enough, but a reliable, efficient, and highly generic implementation is difficult to achieve. We propose the following requirements:

- The start and stop values can be passed as different types, but internally should be of the same type.
- Ranges should work with ordinal types, not just numbers (examples include characters, pointers, and calendar dates).
- If any of the arguments is a floating-point number, a special `FloatRange` type designed to cope well with roundoff is returned.

In the case of ordinal types, the step value is naturally of a different type than the elements of the range. For example, one may add 1 to a character to get the “next” encoded character, but it does not make sense to add two characters.

It turns out that the desired behavior can be achieved with six definitions:

First, given three floats of the same type we can construct a `FloatRange` right away:

```
colon{T<:FloatingPoint}(start::T, step::T, stop::T) = FloatRange{T}(start, step, s
```

Next, if `a` and `b` are of the same type and there are no floats, we can construct a general range:

```
colon{T}(start::T, step, stop::T) = StepRange(start, step, stop)
```

Now there is a problem to fix: if the first and last arguments are of some non-floating-point numeric type, but the step is floating point, we want to promote all arguments to a common floating point type. We must also do this if the first and last arguments are floats, but the step is some other kind of number:

```
colon{T<:Real}(a::T, s::FloatingPoint, b::T) = colon(promote(a,s,b)...)

```

```
colon{T<:FloatingPoint}(a::T, s::Real, b::T) = colon(promote(a,s,b)...)

```

These two definitions are correct, but ambiguous: if the step is a float of a different type than `a` and `b` both definitions are equally applicable. We can add the following disambiguating definition:

```
colon{T<:FloatingPoint}(a::T, s::FloatingPoint, b::T) = colon(promote(a,s,b)...)

```

All of these five definitions require `a` and `b` to be of the same type. If they are not, we must promote just those two arguments, and leave the step alone in case we are dealing with ordinal types:

```
colon{A,B}(a::A, s, b::B) = colon(convert(promote_type(A,B),a), s, convert(promote_type(A,B),b))
```

This example shows that it is not always sufficient to have a built-in set of “promoting operators”. Library functions like this `colon` need more control.

## Current approaches

Numbers tend to be among the most complex features of a language. Numeric types usually need to be a special case: in a typical language with built-in numeric types, describing their behavior is beyond the expressive power of the language itself. For example, in C arithmetic operators like `+` accept multiple types of arguments (ints and floats), but no user-defined C function can do this (this situation is of course improved in C++). In Python, a special arrangement is made for `+` to call either an `__add__` or `__radd__` method, effectively providing double-dispatch for arithmetic in a language that is idiomatically single-dispatch.



### 4.3.3 Multidimensional array indexing

One-dimensional arrays are a simple and essential data structure found in most programming languages. The multi-dimensional arrays required in scientific computing, however, are a different beast entirely. Allowing any number of dimensions entails a significant increase in complexity. Why? The essential reason is that core properties of the data structure no longer fit in a constant amount of space. The space needed to store the sizes of the dimensions (the array shape) is proportional to the number of dimensions. This does not seem so bad, but becomes a large problem due to three additional facts.

First, code that operates on the dimension sizes needs to be highly efficient. Typically the overhead of a loop is unacceptable, and such code needs to be fully unrolled. Second, in some code the number of dimensions is a *dynamic* property — it is only known at run time. Third, programs may wish to treat arrays with different numbers of dimensions very differently. A vector (1d) might have rather different behaviors than a matrix (2d) (for example, to compute a norm). This kind of behavior makes the number of dimensions a crucial part of program semantics, preventing it from remaining a compiler implementation detail.

These facts pull in different directions. The first fact asks for static analysis. The second fact asks for run-time flexibility. The third fact asks for dimensionality to be part of the type system, but partly determined at run time (for example, via virtual method dispatch). Current approaches choose a compromise. In some systems, the number of dimensions has a strict limit (e.g. 3 or 4), so that separate classes for each case may be written out in full. Other systems choose flexibility, and just accept that most or all operations will be dynamically dispatched. Other systems might provide flexibility only at compile time, for example a template library where the number of dimensions must be statically known.

Whatever decision is made, rules must be defined for how various operators act on dimensions. For now we will focus on indexing, since selecting parts of arrays has particularly rich behavior with respect to dimensionality. For example, if a single row

or column of a matrix is selected, does the result have one or two dimensions? Array implementations prefer to invoke general rules to answer such questions. Such a rule might say “dimensions indexed with scalars are dropped”, or “trailing dimensions of size one are dropped”, or “the rank of the result is the sum of the ranks of the indexes” (as in APL).

Our goal here is a bit unusual: we are not concerned with which rules might work best, but merely with how they can be specified, so that domain experts can experiment.

In fact different domains want different things. E.g. in images, each dimension might be quite different, e.g. time vs. space vs. color, so you don’t want to drop or rearrange dimensions very often.

Here are our ground rules:

1. You can’t manually implement the behavior inside the compiler
2. The compiler must be able to reasonably understand the program
3. The code must be reasonably easy to write

How are such rules implemented? For a language with built-in multidimensional arrays, the compiler will analyze indexing expressions and determine an answer using hard-coded logic. However, this approach is not satisfying: we would rather implement the behavior in libraries, so that different kinds of arrays may be defined, or so that rules of similar complexity may be defined for other kinds of objects. But these kinds of rules are unusually difficult to implement in libraries. If a library writes out its indexing logic using imperative code, the host language compiler is not likely to be able to analyze it. Using compile-time abstraction (templates) would provide better performance, but such libraries tend to be difficult to write (and read), and the full complement of indexing behavior expected by technical users strains the capabilities of such systems.

Our dispatch mechanism permits a novel solution. If a multiple dispatch system supports variadic functions and argument “splicing” (the ability to pass a structure

of  $n$  values as  $n$  separate arguments to a function), then indexing behavior can be defined as method signatures.

This solution is still a compromise among the factors outlined above, but it is a new compromise that provides a net-better solution.

Below we define a function `index_shape` that computes the shape of a result array given a series of index arguments. We show three versions, each implementing a different rule that users in different domains might want:

```
# drop dimensions indexed with scalars
index_shape() = ()
index_shape(i::Real, I...) = index_shape(I...)
index_shape(i, I...) = tuple(length(i), index_shape(I...)...)

# drop trailing dimensions indexed with scalars
index_shape(i::Real...) = ()
index_shape(i, I...) = tuple(length(i), index_shape(I...)...)

# rank summing (APL)
index_shape() = ()
index_shape(i, I...) = tuple(size(i)..., index_shape(I...)...)
```

Inferring the length of the result of `index_shape` is sufficient to infer the rank of the result array.

These definitions are concise, easy to write, and possible for a compiler to understand fully using straightforward techniques.

Here is a sample derivation for the call `index_shape(1:m,1,1:n)` (the argument type tuple is `(Range1,Int,Range1)`), using the first definition above (dropping scalar-indexed dimensions):

```
index_shape(1:n, 1, 1:m) => tuple(length(::Range1)::Int, index_shape((::Int, ::Range1)...))

index_shape((::Int, ::Range1)...) => index_shape(::Int, ::Range1)
```

```
index_shape(::Int, ::Range1) => index_shape((::Range1),...)
```

```
index_shape((::Range1),...) => index_shape(::Range1)
```

```
index_shape(::Range1) => tuple(length(::Range1)::Int, index_shape(()...)...)
```

```
index_shape(()...) => index_shape()::()
```

```
back substitute => tuple(length(::Range1)::Int, index_shape(()...)::()...)::(Int,)
```

```
back substitute => tuple(length(::Range1)::Int, ::(Int,)...)
```

```
tuple(::Int, ::(Int,)...) => tuple(::Int, ::Int)
```

```
::(Int, Int)
```

The result type is determined using only dataflow type inference, plus a rule for splicing an immediate container (the type of `f((a,b)...)`  is the type of `f(a,b)`). Argument list destructuring takes place inside the type intersection operator used to combine argument types with method signatures.

This approach does not depend on any heuristics. Each call to `index_shape` simply requires one recursive invocation of type inference. This process reaches the base case `()` for these definitions, since each recursive call handles a shorter argument list (for less-well-behaved definitions, we might end up invoking a widening operator instead).

```
diverge() = randbool() ? () : tuple(1, diverge()...)
```

#### **4.3.4 Array views**

#### **4.3.5 Units**

#### **4.3.6 Even more elimination?**

Some features of the language could be even further eliminated. For example data types could be implemented in terms of lambda abstractions. But certain patterns are so useful that they might as well be provided in a standard form. It also probably makes the compiler much more efficient not to need to pass around and repeatedly analyze full representations of the meanings of such ubiquitous constructs.

# Chapter 5

## Case studies

### 5.1 Numerical linear algebra

Multiple dispatch on special matrices

### 5.2 Grid.jl, computational libraries

### 5.3 JUMP

Metaprogramming tools reused for symbolic algebra

### 5.4 ACAS

### 5.5 Parallel computing, e.g. parallel prefix

### 5.6 IJulia

# Chapter 6

## Conclusion

A nominal conclusion

Here is a snapshot of Julia in April, 2014

Show all the stuff that will be fun to look at because it'll be so antiquated by 2019 by following this kind of approach computer scientists can have it both ways, and make new interesting things that are also immediately useful.

A generation of dynamic languages have been designed by trying variants of the class-based object oriented paradigm. This process has been aided by the development of standard techniques (e.g. bytecode VMs) and reusable infrastructure such as code generators, garbage collectors, and whole VMs like the JVM and CLR.

It is possible to envision a future generation of languages that generalize this design to set-theoretic subtyping instead of just classes. This next generation will require its own new tools, such as partial evaluators (already under development in PyPy and Truffle). One can also imagine these future language designers wanting reusable program analyses, and tools for developing lattices and their operators.

It is interesting to observe that the data model of a language like Julia consists of two key relations: the subtype relation, which is relatively well understood and enjoys useful properties like transitivity, but also the typeof relation, which relates individual values to their types (i.e. the ‘typeof’ function). The typeof relation appears not to be transitive, and also has a degree of arbitrariness: a value is of a type merely because it is labeled as such, and because various bits of code conspire to ensure that this

labeling makes sense according to various criteria.

We have speculated about whether future languages will be able to do away with this distinction. One approach is  $\lambda_{aleph}$ . We have also speculated that this could be done using types based on non-well-founded set theory, combining the subset-of and element-of relations using self-containing sets. We are not yet sure what a practical language based on this idea might look like.



# Bibliography

- [1] James Hiram Jr. Morris. *Lambda-Calculus Models of Programming Languages*. PhD thesis, Massachusetts Institute of Technology, 1968.