

Documentation Generation as Information Visualization

Will Crichton*

November 2020
CMU-ISR-20-115E

Institute for Software Research
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213

*This paper was presented at PLATEAU 2020:
The 11th Annual Workshop on the Intersection of HCI and PL
November 2020, Co-located with SPLASH 2020*

*Stanford University

Keywords: documentation generation, information visualization

Abstract

Automatic documentation generation tools, or auto docs, are widely used to visualize information about APIs. However, each auto doc tool comes with its own unique representation of API information. In this paper, I use an information visualization analysis of auto docs to generate potential design principles for improving their usability. Developers use auto docs as a reference by looking up relevant API primitives given partial information, or leads, about its name, type, or behavior. I discuss how auto docs can better support searching and scanning on these leads, e.g. by providing more information-dense visualizations of method signatures.

Documentation Generation as Information Visualization

Will Crichton

Stanford University

wrichto@cs.stanford.edu

Abstract

Automatic documentation generation tools, or auto docs, are widely used to visualize information about APIs. However, each auto doc tool comes with its own unique representation of API information. In this paper, I use an information visualization analysis of auto docs to generate potential design principles for improving their usability. Developers use auto docs as a reference by looking up relevant API primitives given partial information, or leads, about its name, type, or behavior. I discuss how auto docs can better support searching and scanning on these leads, e.g. by providing more information-dense visualizations of method signatures.

2012 ACM Subject Classification Software and its engineering → Documentation; Human-centered computing → Information visualization

Keywords and phrases documentation generation, information visualization

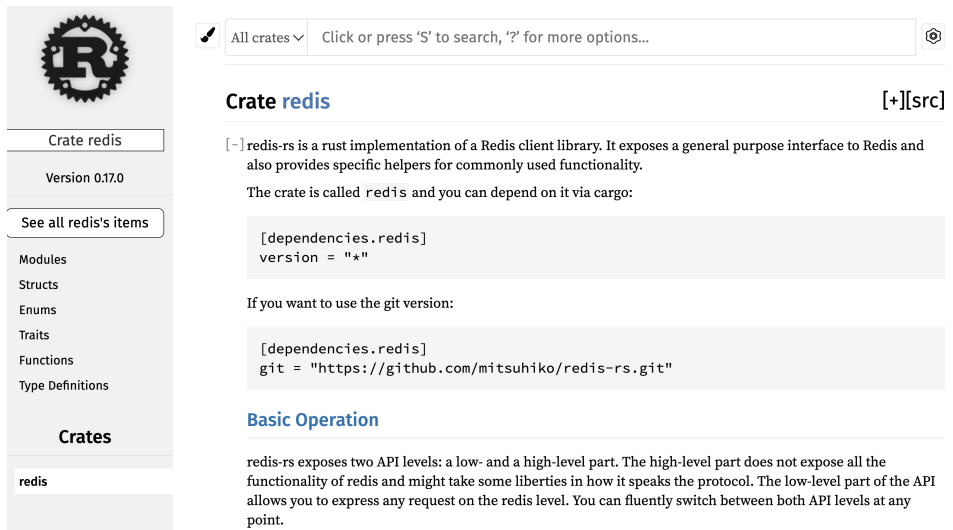
1 Introduction

Understanding other peoples' code is a fact of life in modern software development. With the rise of centralized package repositories, it's easier than ever to access third party libraries and frameworks. In 2019, the average JavaScript package had five direct dependencies and 86 transitive dependencies [6]. Developers cite the availability of libraries as the #1 reason to adopt a programming language [4]. Hence, making libraries easier to understand is a chief usability concern in today's programming landscape.

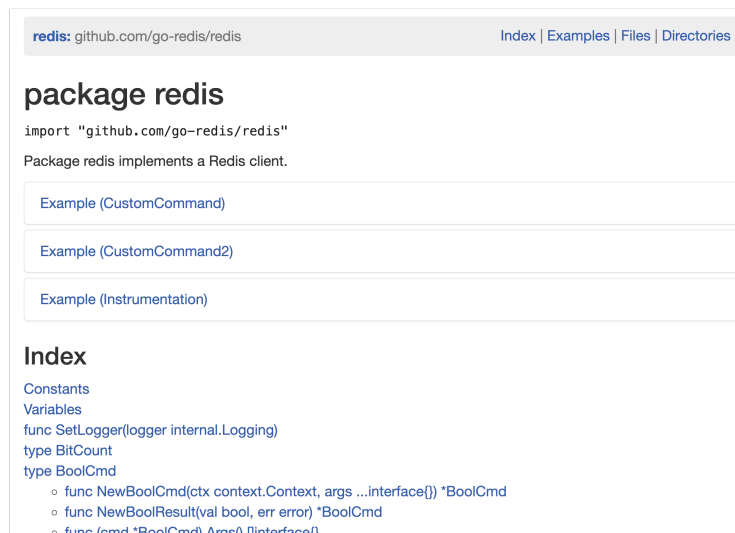
For the vast majority of libraries, their clients will never actually read the library source code. Instead, programmers rely on external resources: documentation, examples, tutorials, StackOverflow, blog posts, coworkers, and so on. However, these resources all have the same problem: they require a human. Someone has to write the blog, answer the StackOverflow post, and carefully craft the doc comments. Worse, that person probably has no training in technical communication, information visualization, computing education, or any discipline relevant to effectively explaining software. And most problematically, these resources can become out of date as soon as the library changes.

Automatically generated documentation, or auto docs, solve these problems by being derived directly from the source code. Tools such as Javadoc or Rustdoc take a codebase as input, and produce a wiki-like website as output. Except for doc comments, auto docs are always up-to-date. The author of the auto doc software can (potentially) present information according to best practices. And all this comes for free — no extra effort needed by the library author! Every major programming language has a documentation generator, and recent languages consider one so vital that they ship it with the official toolset. Auto docs are arguably the most successful and widely-used software visualization tool in history.

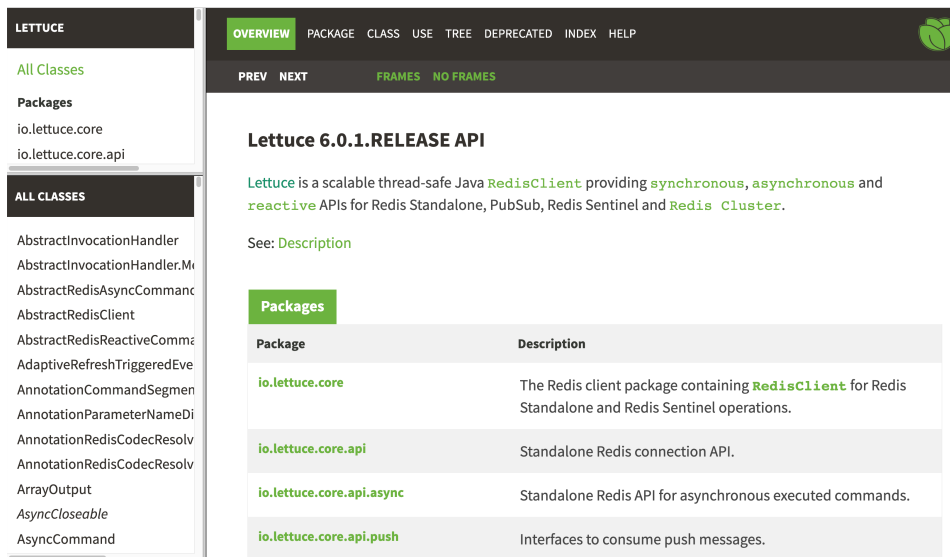
Despite their apparent importance, auto docs are rarely the subject of research. They have few established best practices. Each new programming language comes with an entirely new docs format, as shown in Figure 1. Beyond aesthetics, each interface provides significantly different means of visualizing and searching information. This variation raises the question: is one format better than another? In this paper, I show how an information visualization analysis can generate potential design principles for auto docs.



The screenshot shows the crates.io page for the `redis` crate. On the left is a sidebar with navigation options like 'Crate redis', 'Version 0.17.0', and 'See all redis's items'. The main content area includes a search bar, the crate name 'Crate redis', a description of the crate as a Rust implementation of a Redis client library, and dependency information. A 'Basic Operation' section explains the two API levels (low and high) and how to switch between them.



The screenshot shows the Go documentation page for the `redis` package. It features a navigation bar with links to 'Index', 'Examples', 'Files', and 'Directories'. The main content includes the package name 'package redis', an import statement, a description of the package as a Redis client, and several example code snippets. An 'Index' section lists constants, variables, and functions like `SetLogger`, `BitCount`, and `BoolCmd`.



The screenshot shows the Java documentation page for the 'Lettuce' library. It has a navigation bar with tabs for 'OVERVIEW', 'PACKAGE', 'CLASS', 'USE', 'TREE', 'DEPRECATED', 'INDEX', and 'HELP'. The main content includes the title 'Lettuce 6.0.1.RELEASE API', a description of the library as a scalable thread-safe Java Redis client, and a 'Packages' table. The table lists packages like `io.lettuce.core`, `io.lettuce.core.api`, `io.lettuce.core.api.async`, and `io.lettuce.core.api.push` with their respective descriptions.

■ **Figure 1** A screenshot of the front page of automatically-generated documentation for a Redis library in different languages. From top to bottom: Rust, Go, Java.

2 Task Analysis

To analyze auto docs from an information visualization perspective, we need to understand three things: the information, the tasks, and the representations that bridge the two. Then we can identify the gap between theory and reality to generate design ideas for future auto doc tools.

2.1 What information is being represented?

Auto docs generally present information about APIs, i.e. aspects that are externally visible to consumers of the library. APIs generally include:

- **Data structures:** structs, enums, classes
- **Functions:** constructors, methods, standalone functions
- **Interfaces:** abstract classes, traits, typeclasses
- **Hierarchies:** modules, files, namespaces

These code constructs collectively describe the primitives of an API. A developer can usually provide source-code comments further explaining their purpose, which may contain prose description or code examples. Additionally, many relations arise from the structure of these objects. For example:

- **Inputs:** a type is used as input to a function
- **Outputs:** a type is used as output of a function
- **Contains:** a type is a field of another type
- **Inherits:** a class inherits from another class
- **Implements:** a type implements an interface

Besides the explicit hierarchies within the language (e.g. namespaces), relations implicitly define hierarchies amongst API primitives. For example, the C++ auto docs tool Doxygen will generate an inheritance tree diagram of all classes in a codebase. Rust's auto docs will also show all the types that implement a trait on the trait's generated page.

2.2 What tasks use this information?

Auto docs work best as a reference. That is to say, auto docs provide detailed information about individual constructs within an API. By contrast, tutorials are more effective at showing how multiple API constructs can be combined to accomplish a larger task. Therefore we will focus on tasks that developers have specifically for reference information.

The key feature of a reference is its ability to guide a reader with partial information about their object of interest. For example, if a person remembers the first letter of a particular word in a textbook, they can use a glossary to find the relevant word and corresponding pages. More generally, a person starts with a lead (e.g. a letter), and a reference provides some level of guidance in converting the lead into more information (e.g. a table of words grouped by first letter). For auto docs, the question is then: what leads do developers use when seeking information? Here are some example scenarios:

- A developer wants to turn a list into a set. Their lead is *functions that return the Set type and possibly take the List type as input*.
- A developer wants to find an element of a list that matches a predicate. Their lead is *the natural language keywords "list", "find", and "predicate"*.

- A developer wants to know what immutable operations exist for lists. Their lead is *methods on the List type that have a read-only modifier*.

Information foraging analyses have shown that programmers use leads to search documentation [3]. However, no prior work has compiled a holistic taxonomy of such leads — an appropriate starting point for future work on auto docs.

2.3 What is the best information representation for each task?

Consider a developer with a lead of natural language keywords, e.g. the “find” and “predicate” example. They might adopt a range of strategies to forage for this information, for example in this order:

1. **Search engine:** they go to Google and type “<language> list find predicate”. No relevant results show up.
2. **Browser search:** they go to the auto doc page for the List data structure. They Ctrl+F in the browser for each keyword, but don’t find a relevant method.
3. **Scanning:** on the auto doc page, they scan through the list of method names and description. Realizing the method was called “indexOf” instead of “find”, they locate the appropriate method.


These examples highlights two key needs of a developer following a lead. First, they need to be able to encode the lead into a search. Second, if the search fails, they need a visualization of many possibly relevant objects which can be manually searched. Then the developer can make fuzzier associations between their lead and the provided information, e.g. the semantic relationship between “indexOf” and “find”.

2.4 Where do current information representations fall short?

For search, most auto docs tools support (at best) natural language keywords. But many relational leads are difficult to encode in a query. “Functions that take a list as input and return an integer” has a formal representation: the data type `[a] -> int` where `[a]` means a list of any kind of element. Notably, Haskell has a type-based search engine, Hoogle [1], which enables these kinds of queries. But no other language or auto doc tool supports these queries.

For scanning, auto docs tools have two issues: first, scanning requires an initial filter or anchor to limit the set of possible matches. For most tools today, functions are anchored around their class. That means a developer can easily get a webpage of all the methods on the List class, but as mentioned above, they cannot easily get a page of all methods that return a List. Even within a class’s method list, it can be valuable to filter. For example, Rust’s documentation for `Vec<T>` [2] lists well over 100 methods. The docs provides no way to e.g. filter for read-only methods that take `&self` as input.

Second, auto doc tools do not necessarily provide efficient information representations to facilitate rapid scanning. Most auto docs today present methods in a one-dimensional list. However, a developer may more easily scan a large number of method names if they can all be on screen at once, e.g. in a two-dimensional table like the prototype in Figure 2. Alternatively, organizing methods hierarchically could turn a linear search into logarithmic. For example, the methods of a stateful class could be grouped around which states they apply to, which Sunshine et al. [5] demonstrated will empirically reduce documentation search times for certain tasks.



&self	&mut self	static
capacity(&Self) -> usize	reserve(&mut Self, usize) -> Result<(), TryReserveError>	new() -> Vec<T>
as_slice(&Self) -> &[T]	reserve_exact(&mut Self, usize) -> Result<(), TryReserveError>	with_capacity(usize) -> Vec<T>
as_ptr(&Self) -> *const T	shrink_to_fit(&mut Self) -> bool	from_raw_parts(*mut T, usize, usize) -> Vec<T>
len(&Self) -> usize	shrink_to(&mut Self, usize) -> bool	leak(Vec<T>) -> &'a mut [T]
is_empty(&Self) -> bool	truncate(&mut Self, usize) -> bool	from(T) -> T
borrow(&Self) -> &T	as_mut_slice(&mut Self) -> &mut [T]	try_from(U) -> Result<T, <T as TryFrom<U>>::Error>
type_id(&Self) -> TypeId	as_mut_ptr(&mut Self) -> *mut T	from(&'a Vec<T>) -> Cow<'a, [T]>
to_owned(&Self) -> T	as_ref(&Self) -> &Vec<T>	from(String) -> Vec<u8>
clone_into(&Self, &mut T) -> Option<Ordering>	as_ref(&Self) -> &[T]	from(Vec<T>) -> Cow<'a, [T]>
partial_cmp(&Self, &Vec<T>) -> Option<Ordering>	set_len(&mut Self, usize) -> bool	from(Cow<'a, [T]>) -> Vec<T>
index(&Self, I) -> &Vec<T> as Index<I>::Output	swap_remove(&mut Self, usize) -> T	from([T; N]) -> Vec<T>
fmt(&Self, &mut Formatter) -> Result<(), Error>	insert(&mut Self, usize, T) -> bool	from(VecDeque<T>) -> Vec<T>
cmp(&Self, &Vec<T>) -> Ordering	remove(&mut Self, usize) -> T	from(&str) -> Vec<u8>
as_ref(&Self) -> &Vec<T>	retain(&mut Self, F) -> bool	from(BinaryHeap<T>) -> Vec<T>
as_ref(&Self) -> &[T]	dedup_by_key(&mut Self, F) -> bool	from(Vec<T>) -> Rc<[T]>
deref(&Self) -> &[T]	dedup_by(&mut Self, F) -> bool	from(Box<[T]>) -> Vec<T>
hash(&Self, &mut H) -> bool	push(&mut Self, T) -> Option<T>	from(Vec<T>) -> Arc<[T]>
borrow(&Self) -> &[T]	pop(&mut Self) -> Option<T>	from(&mut [T]) -> Vec<T>
eq(&Self, &Vec) -> bool	append(&mut Self, &mut Vec<T>) -> bool	from(Vec<T>) -> VecDeque<T>
ne(&Self, &Vec) -> bool	drain(&mut Self, R) -> Drain<T>	from(Vec<T>) -> Box<[T]>
eq(&Self, &Vec) -> bool	clear(&mut Self) -> bool	from_iter(I) -> Vec<T>
ne(&Self, &Vec) -> bool	split_off(&mut Self, -> Vec<T>	default() -> Vec<T>
eq(&Self, &Vec) -> bool		from(CString) -> Vec<u8>
ne(&Self, &Vec) -> bool		from(Vec<NonZeroU8>) -> CString
eq(&Self, &Vec) -> bool		
ne(&Self, &Vec) -> bool		

Figure 2 Prototype of a scanning-oriented interface for Rust’s auto docs. By displaying methods in a table and grouping by the first argument, developers can quickly search for relevant methods.

3 Conclusion

My goal in this paper is to provide an initial framework for generating ideas about improving auto docs. Like many programming tools, auto docs are widely used, but have yet to be examined critically from a human-centered perspective. As Bret Victor points out in “Learnable Programming” [7], an API resource should “dump the parts bucket onto the floor.” However, developers shouldn’t have to spend hours searching through a sea of parts. With a careful understanding of what developers are looking for and what leads they use, we can build better tools for searching and visualizing API information.

References

- 1 Hoogle, 2020. URL: <https://hoogle.haskell.org/>.
- 2 std::vec::Vec - Rust, 2020. URL: <https://doc.rust-lang.org/std/vec/struct.Vec.html>.
- 3 Amy J Ko, Brad A Myers, Michael J Coblenz, and Htet Htet Aung. An exploratory study of how developers seek, relate, and collect relevant information during software maintenance tasks. *IEEE Transactions on software engineering*, 32(12):971–987, 2006.
- 4 Leo A Meyerovich and Ariel S Rabkin. Empirical analysis of programming language adoption. In *Proceedings of the 2013 ACM SIGPLAN international conference on Object oriented programming systems languages & applications*, pages 1–18, 2013.
- 5 Joshua Sunshine, James D Herbsleb, and Jonathan Aldrich. Structuring documentation to support state search: A laboratory experiment about protocol programming. In *European Conference on Object-Oriented Programming*, pages 157–181. Springer, 2014.
- 6 Ruturaj K Vaidya, Lorenzo De Carli, Drew Davidson, and Vaibhav Rastogi. Security issues in language-based software ecosystems. *arXiv preprint arXiv:1903.02613*, 2019.
- 7 Bret Victor. Learnable programming, 2012. URL: <http://worrydream.com/LearnableProgramming/>.