# Tools for Disambiguating RFCs

Jane Yen, Barath Raghavan, Ramesh Govindan

## **USC**Viterbi

School of Engineering

### **RFC** production

YEAR	2016	2017	2018	2019	2020	2021
RFC COUNT	310	263	208	180	209	240

### Specification production

We need to standardize this XYZ protocol.

XYZ protocol uses three way handshake...



Communication



Working Group

**Specification Author** 

### Ambiguities



Ambiguous Specification Human Interpretation

Multiple versions

### **Classic examples**

The checksum is the 16-bit one's complement of the one's complement sum of the ICMP message starting with the ICMP Type.

### **Classic examples**

The checksum is the 16-bit one's complement of the one's complement sum of the ICMP message <u>starting</u> with the ICMP Type.

Ending at ?

### **Classic examples**

The checksum is the 16-bit one's complement of the one's complement sum of the ICMP message <u>starting</u> with the ICMP Type.

OR

#### Ending at ?





Checksum both the header and payload

### Possible Consequences



**Buggy Implementation** 



### **Rigorous discussion**

This sentence is ambiguous. Please rephrase.

The \$BOO field MUST be 0



Communication \* N



Working Group

**Specification Author** 

### Are we close to near 0-ambiguity specification?

### Our work

#### Semi-Automated Protocol Disambiguation and Code Generation Tamás Lévai

Jane Yen University of Southern California yeny@usc.edu

and Economics levait@tmit.bme.hu

Xiang Ren University of Southern California xiangren@usc.edu

#### ABSTRACT

For decades, Internet protocols have been specified using natural language. Given the amhiguity inherent in such text, it is not our prising that protocol implementations have long exhibited bugs In this paper, we apply natural language processing (NLP) to el fect semi-automated generation of protocol implementations from specification text. Our system, SAGE, can uncover ambiguous or under-specified sentences in specifications; once these are clari fied by the author of the protocol specification, SAGE can generate protocol code automatically.

of under-specification in the ICMP RFC; after fixing these, sada is able to automatically generate code that interoperates perfectly with Linux implementations. We show that sace generalizes to sections of BFD, IGMP, and NTP and identify additional conceptual components that SAGE needs to support to generalize to complete complex protocols like BGP and TCP.

CCS CONCEPTS

Networks → Formal specifications;

#### KEYWORDS

natural language, protocol specifications

ACM Reference Formati Jane Yen, Tamis Lévai, Qinyuan Ye, Xiang Ren, Ramesh Govindar and Barath Raghavan. 2021. Semi-Antemated Pretocol Disambiguation and Code Generation. In ACM SIGCOMM 2021 Conference (SIGCOMM '21) Aurust 23-28, 2021, Virtual Event, USA, ACM, New York, NY, USA, 15 pares https://doi.org/10.1145/3452296.3472910

#### 1 INTRODUCTION

Four decades of Internet protocols have been specified in English and used to create, in Clark's words, rough consensus and running code [16]. In that time we have come to depend far more on network protocols than most imagined. To this day, engineers implement a protocol by reading and interpreting specifications as described in Resuest For Comments documents (RFCs). Their challenge is



SECOMM '22, August 23-28, 2022, Virtual Event, USA 021 Copyright held by the owner/au M ISBN 978-1-4503-8383-7/21/38. er//doi.org/10.1145/3452296.3472910

Qinyuan Ye Budapest University of Technology Uni versity of Southern California qinyuany@usc.edu

Software engineers find it difficult to interpret specifications in

large part because natural language can be ambiguous. Unfortu-

nately, such ambiguity is not rate: the errata alone for RFCs over

have caused [17, 33, 68, 78]. Ambiguity has resulted in buggy imple-

mentations, security vulnerabilities, and has necessitated expensiv

and time-consuming software engineering processes, like interop-

protocol implementations from RFCs. Our main challenge is to

understand the semantics of a specification. This task semantic

pursing, has advanced in recent years with parsing tools such as

CCG [5]. Such tools describe natural language with a lexicon and

yield a semantic interpretation for each sentence. Because they

are trained on generic prose, they cannot be expected to work out

of the box for idiomatic network protocol specifications, which

fields), incomplete sentences, and implicit context from neighboring

text or other protocols. More importantly, the richness of natural

language will likely always lead to ambiguity, so we do not expect

Contributions. In this paper, we describe SAGE, a semi-automated

either case, the user (e.g., the author of the specification) can then

revise the sentences and re-run sada until the resulting REC can

cleanly be turned into code. SAGE can be used at various stages

fully-automated NLP-based systems (\$2)

contain embedded syntactic cues (e.g., structured descriptions of

Ramesh Govindan **Barath Raghavan** University of Southern California University of Southern California ramesh@usc.edu barathra@usc.edu to navigate easy-to-misinterpret colloquial language while writing not only a bug-free implementation but also one that interoperates

with code written by another person at a different time and place. the years highlight numerous ambiguities and the problems they Using SAGE, we discover 5 instances of ambiguity and 6 instances

erability bake-offs [32, 71]. To address this, one line of research has sought formal specification of programs and protocols (§8), which would enable verifying specification correctness and, potentially, enable automated code generation [13]. However, formal specifications are cumbersome and thus have not been adopted in practice; to date, protocols are specified in natural language.<sup>1</sup> In this paper, we apply NLP to semi-automated generation of

approach to protocol analysis and code generation from naturallanguage specifications. SAGE reads the natural-language protocol specification (e.g., an RFC or Internet Draft) and marks sentences (a) for which it cannot generate unique semantic interpretations or (b) which fail on the protocol's unit tests (SAGE uses test-driven development). The former sentences are likely semantically ambiguous whereas the latter represent under-specified behaviors. In



<sup>1</sup>In recent years, attempts have been made to formalize other aspects of network operation, such as network configuration [7, 37] and control plane behavior [57], with

272

SIGCOMM'21

#### venv@usc.edu Abstract

For decades, drafting Internet protocols has taken significant amounts of human supervision due to the fundamental ambiguity of natural language. Given such ambiguity, it is also not surprising that protocol implementations have long exhibited bugs. This pain and overhead can be significantly reduced with the help of natural language processing (NLP). We recently applied NLP to identify ambiguous or underspecified sentences in RFCs, and to generate protocol implementations automatically when the ambiguity is clarified. However this system is far from general or deployable. To further reduce the overhead and errors due to ambiguous sentences, and to improve the generality of this system, much work remains to be done. In this paper, we consider what it would take to produce a fully-general and useful system for easing the natural-language challenges in the RFC process.

#### **CCS** Concepts

• Networks  $\rightarrow$  Formal specifications. Keywords

Jane Yen

#### natural language, protocol specifications

ACM Reference Format: Jane Yen Ramesh Gowindan and Barath Rashavan 2021 Tools for Disambiguating RFCs. In Applied Networking Research Workshop (ANRW '21), July 24-30, 2021, Virtual Event, USA. ACM, New York, NY, USA, 7 pages. https://doi.org/10.1145/3472305.3472314

#### 1 Introduction

It has long been the case that protocol behaviors are discussed and debated in the networking community for years before they are codified. The Internet Architecture Board (IAB) [9], Internet Engineering Task Force (IETF), and Internet Research Task Force (IRTF) [12] enable groups of independent networking professionals to draft Request for



This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike International 4.0 License ANRW '21. July 24-30. 2021. Virtual Event, USA @ 2021 Copyright held by the owner/author(s ACM ISBN 978-1-4503-8618-0/21/07 https://doi.org/10.1145/3472305.3472314

Ramesh Govindan Barath Raghavan University of Southern California University of Southern California University of Southern California

**Tools for Disambiguating RFCs** 

ramesh@usc.edu barathra@usc.edu

> Comments (RFCs) to explain their networking ideas in English. Their aim is to publish their ideas globally and, typically to codify those ideas as a networking standard. Draft RFCs are carefully evaluated and reviewed by members of working groups and editors; these individuals manually consider editorial and technical perspectives. Though the process by which RFCs are developed is effective, and can identify significant vague, incorrect, or partial descriptions and get authors to their drafts, the overall overhead of this human supervision is significant.

> RFC editors follow a number of guidelines to review RFC drafts and provide feedback to the authors; authors typically aim to satisfy the guidelines as guickly as possible so that the back-and-forth review process can be minimized. The guidelines mostly address what sections are required to be present, what order the sections should be in, what minimal infor mation should be stated clearly, and so forth. While these guidelines help drafts across time and space to maintain a uniform style, these guidelines are mostly writing advice and do not effectively help authors or reviewers to identify technical issues. To tackle technical concerns in any RFC draft requires participants' professional background knowledge RFC authors have many ways to describe technical details including natural language, pseudocode, formal specifications, real code, state machine diagrams, and more, While pseudocode, formal specifications, and real code can provide precise execution steps, natural language is still preferable due to its flexibility and readability. However, natural lan guage can also be fuzzy, which sometimes leads to a fun damental technical problem. For example, "the type code changed to 0" is a verbatim sentence in Internet Control Message Protocol (ICMP) RFC [23]. The noun phrase "type code" can confuse a reader as to whether it refers to as a type named "code" or a code named "type" or a specific protocol header field named "type code". If any reader interprets the noun phrase incorrectly, a packet may be generated incor rectly and get dropped by a receiver.

> To strike the right balance when using natural language, we might ask the following question: Can we systematically identify natural language ambiguity in network protocol specifications and make changes accordingly? The answer to this question would help many stakeholders in the context of protocol specifications For example, working groups and standards writers could leverage techniques to avoid or reduce back-and-forth

**ANRW'21** 

85

### Sage



Oinvuan Ye

University of Southern California

qinyuany@usc.edu

#### Semi-Automated Protocol Disambiguation and Code Generation

ramesh@usc.edu

272

SIGCOMM'21

Jane Yen	Tamás Lévai
University of Southern California	Budapest University of Technology
yeny@usc.edu	and Economics
	levait@tmit.bme.hu

Xiang Ren University of Southern California xiangren@usc.edu

#### ABSTRACT

For decades, Internet protocols have been specified using natural language. Given the ambiguity inherent in such text, it is not surprising that protocol implementations have long exhibited bugs In this paper, we apply natural language processing (NLP) to el fect semi-automated generation of protocol implementations from specification text. Our system, SAGE, can uncover ambiguous or under-specified sentences in specifications; once these are clari fied by the author of the protocol specification, SAGE can generate protocol code automatically. Using SAGE, we discover 5 instances of ambiguity and 6 instances

of under-specification in the ICMP RFC; after fixing these, sage is able to automatically generate code that interoperates perfectly with Linux implementations. We show that sage generalizes to sections of BFD, IGMP, and NTP and identify additional conceptual components that SAGE needs to support to generalize to complete complex protocols like BGP and TCP.

CCS CONCEPTS

Networks → Formal specifications;

#### KEYWORDS

natural language, protocol specifications

ACM Reference Formati Jane Yen, Tamis Lévai, Qinyuan Ye, Xiang Ren, Ramesh Govindar and Barath Raghavan. 2021. Semi-Antemated Pretocol Disambiguation and Code Generation. In ACM SIGCOMM 2021 Conference (SIGCOMM '21). Aurust 23-28, 2021, Virtual Event, USA, ACM, New York, NY, USA, 15 pares https://doi.org/10.1145/3452296.3472910

#### 1 INTRODUCTION

Four decades of Internet protocols have been specified in English and used to create, in Clark's words, rough consensus and running code [16]. In that time we have come to depend far more on network protocols than most imagined. To this day, engineers implement a protocol by reading and interpreting specifications as described in Resuest For Comments documents (RFCs). Their challenge is



SIGCOMM '22, August 23-28, 2022, Virtual Event, USA © 2021 Copyright held by the owner and ACM ISBN 978-1-4503-4383-7/21/08. https://doi.org/10.1145/3451296.3477910

it@tmit.bme.hu Ramesh Govindan **Barath Raghavan** University of Southern California University of Southern California barathra@usc.edu to navigate easy-to-misinterpret colloquial language while writing not only a bug-free implementation but also one that interoperates with code written by another person at a different time and place. Software engineers find it difficult to interpret specifications in large part because natural language can be ambiguous. Unfortu

nately, such ambiguity is not rate: the errata alone for RFCs over

the years highlight numerous ambiguities and the problems they have caused [17, 33, 68, 78]. Ambiguity has resulted in buggy implementations, security valnerabilities, and has necessitated expensive and time-consuming software engineering processes, like interoperability bake-offs [32, 71]. To address this, one line of research has sought formal specification of programs and protocols (§8), which would enable verifying specification correctness and, potentially, enable automated code generation [13]. However, formal specifications are cumbersome and thus have not been adopted in practice; to date, protocols are

specified in natural language. In this paper, we apply NLP to semi-automated generation of protocol implementations from RFCs. Our main challenge is to understand the semantics of a specification. This task semantic parsing, has advanced in recent years with parsing tools such as CCG [5]. Such tools describe natural language with a lexicon and yield a semantic interpretation for each sentence. Because they are trained on generic prose, they cannot be expected to work out of the box for idiomatic network protocol specifications, which contain embedded syntactic cues (e.g., structured descriptions of fields), incomplete sentences, and implicit context from neighboring text or other protocols. More importantly, the richness of natural language will likely always lead to ambiguity, so we do not expect fully-automated NLP-based systems (\$2) Contributions. In this paper, we describe SAGE, a semi-automated

#### approach to protocol analysis and code generation from naturallanguage specifications. SAGE reads the natural-language protocol specification (e.g., an RFC or Internet Draft) and marks sentences (a) for which it cannot generate unique semantic interpretations or (b) which fail on the protocol's unit tests (SAGE uses test-driven development). The former sentences are likely semantically am biguous whereas the latter represent under-specified behaviors. In either case, the user (e.g., the author of the specification) can then revise the sentences and re-run sada until the resulting REC can

cleanly be turned into code. SAGE can be used at various stages

<sup>1</sup>In recent years, attempts have been made to formalize other aspects of network operation, such as network configuration [7, 37] and control plane behavior [53], with earsing degrees of success.

- Uncover 5 instances of ambiguity and 6 instances of under-specification in ICMP RFC
- Generate executable code of unambiguous specification that interoperate with 3rd party code
- Generalize to sections of BFD, IGMP and NTP

### Our Approach: Use Natural Language Processing

#### Natural language processing (NLP) on English specifications



Specification

Ambiguity Discovery

### Goal

#### Executable protocol code generation



Specification

Ambiguity Discovery

**Executable Code** 

### Human in the Loop

#### Semantic parsing is not perfect



Specification

Ambiguity Discovery

Executable Code

16

• Specifications use domain-specific language

• Specifications use domain-specific language

• Semantic parsers have limitations

• Specifications use domain-specific language

• Semantic parsers have limitations

• Semantic representations need to be converted into code

### Contributions

• Specifications use domain-specific language

Extend semantic parser with domain-specific syntax and semantics

• Semantic parsers have limitations

• Semantic representations need to be converted into code

### Contributions

• Specifications use domain-specific language

Extend semantic parser with domain-specific syntax and semantics

• Semantic parsers have limitations

Automate disambiguation of poor semantic representations with checking rules

• Semantic representations need to be converted into code

### Contributions

• Specifications use domain-specific language

Extend semantic parser with domain-specific syntax and semantics

• Semantic parsers have limitations

Automate disambiguation of poor semantic representations with checking rules

• Semantic representations need to be converted into code

Compile semantic representations into executable code

### Sage Components









### Sage Workflow



### Sage Workflow



## Sage Components



### Semantic parsing



## Key Observation

A logical form is a unifying abstraction for disambiguation and code generation

## **Domain Specific Extensions**

Term dictionary with generic noun or noun phrase labeler

•

Domain specific semantics

## **Domain Specific Extensions**

Term dictionary with generic noun or noun phrase labeler

- Part of speech tagging: SpaCy
- Extended SpaCy's term dictionary
  - e.g., "one's complement"

Domain specific semantics

## **Domain Specific Extensions**

Term dictionary with generic noun or noun phrase labeler

- Part of speech tagging: SpaCy
- Extended SpaCy's term dictionary
  - e.g., "one's complement"

Domain specific semantics

• Idiomatic usage

• e.g., "=" sign in "0 = Echo Reply"

## Sage Components

![](_page_32_Figure_1.jpeg)

## Ambiguity

CCG parser could generate zero or more than one logical forms (LFs)

### Ambiguity

CCG parser could generate zero or more than one logical forms (LFs)

Incomplete

If code = 0, identifies the octet where an error was detected

0 LF

## Ambiguity

CCG parser could generate zero or more than one logical forms (LFs)

Incomplete

If code = 0, identifies the octet where an error was detected

0 LF

Imprecise language To form a information reply message, the source and destination addresses are simply reversed, the type code changed to 16, and the checksum recomputed

code? type?

## Winnowing Ambiguous Logical Forms

![](_page_36_Figure_1.jpeg)

Checking rules

### SAGE Components

![](_page_37_Figure_1.jpeg)

![](_page_38_Figure_0.jpeg)

Logical Forms to Code

### Evaluation

TEST COMMANDS	PURPOSE
client ping -c 10 10.0.1.1	Test echo msg
client ping -c 10 192.168.3.1	Test dest unreachable msg
client ping -c 10 -t 1 192.168.2.2	Test time exceeded msg
client traceroute 10.0.1.1	Test traceroute

![](_page_39_Figure_2.jpeg)

#### **Tools for Disambiguating RFCs**

# Beyond Sage, there remains many challenges unaddressed.

#### **Tools for Disambiguating RFCs**

Jane Yen Ramesh Govindan University of Southern California venv@usc.edu university of Southern California

#### Abstract

For decades, during Internet protocols has taken signifcant amounts of human supervision due to the fundamental ambiguty of natural language. Given such ambiguty, it is also not surprised that postoos implementations have long exhibited bags. This pain and overhead can be significantly releved with the human bags and the second protocols implementations automatically when the ambiguty is clarificat theorem in the second second second second second However this systems in far from general or deloyable. To further relace the overhead and errors due to ambigue sentences, and to improve the generality of this system, much woold take to produce a full-general and useful system for soming the natural language challenges in the FCP processoming the natural language challenges in the FCP process.

#### CCS Concepts

 $\bullet \ Networks \rightarrow Formal \ specifications.$ 

#### Keywords

natural language, protocol specifications

ACM Reference Format: Jane Yen, Rameh Govindan, and Barath Raghavan. 2021. Tools for Disambiguating RFCs. In Applied Networking Research Workshop (ANRW '21), July 24-30, 2021. Virtual Event, USA. ACM, New York, NY, USA, 7 pages. https://doi.org/10.1145/3972305.3472315

#### 1 Introduction

It has long been the case that protocol behaviors are discussed and debated in the networking community for years before they are codified. The literest Architecture Board (IAB) [9], Internet Engineering Task Force (IETP) (21) enable groups of independent networking professionals to draft Request for

![](_page_40_Picture_13.jpeg)

This work is licensed under a Creative Commons Attribution NonCommercial-ShareAlike International 4.0 License. ANRW '21, July 24–30, 2021, Virtual Event, ISA 2021 Copyright held by the owner/author(s). ACM ISBN 976-1-450-4618-462107. Interview over 2010 1105 9427384 9427314

Ramesh Govindan Barath Raghavan University of Southern California massh@usc.edu Comments (RFCs) to explain their netwerking ideas in Eng-

Inh. There aim is to publish their ideas globally and, typically, to codify those ideas as a networking standard: Draft RFGs are carefully evaluated and reviewed by members of working groups and editor; these individuals manually consider editorial and technical perspectives. Though the process significant vagues, incorrect, or partial descriptions and get authors to their drafts, the overall overhead of this human supervision is significant.

RFC editors follow a number of guidelines to review RFC drafts and provide feedback to the authors: authors typically aim to satisfy the guidelines as guickly as possible so that the back-and-forth review process can be minimized. The guide lines mostly address what sections are required to be present. what order the sections should be in, what minimal information should be stated clearly, and so forth. While these guidelines help drafts across time and space to maintain a uniform style, these guidelines are mostly writing advice and do not effectively help authors or reviewers to identify technical issues. To tackle technical concerns in any RFC draft requires participants' professional background knowledge RFC authors have many ways to describe technical details, including natural language, pseudocode, formal specifica tions, real code, state machine diagrams, and more. While pseudocode, formal specifications, and real code can provide precise execution steps, natural language is still preferable due to its flexibility and readability. However, natural language can also be fuzzy, which sometimes leads to a fundamental technical problem. For example, "the type code changed to 0" is a verbatim sentence in Internet Control Message Protocol (ICMP) RFC [23]. The noun phrase "type code" can confuse a reader as to whether it refers to as a type named "code" or a code named "type" or a specific protocol header field named "type code". If any reader interprets the noun phrase incorrectly, a packet may be generated incorrectly and get dropped by a receiver.

### **SAGE** Limitations

• Specification components

Name	Description
Packet Format	Packet anatomy (i.e., field structure)
Field Descriptions	Packet header field descriptions
Constraints	Constraints on field values
Protocol Behaviors	Reactions to external/internal events
System Architecture	Protocol implementation components
+ State Management	Session information and/or status
Comm. Patterns	Message sequences (e.g., handshakes)

Table 1: Protocol specification components. SAGE supports those marked with (fully) and + (partially).

- Paragraph or sentence-based analysis
- Mis-matched/mis-captured behaviors
- Standalone or multiple RFCs
- Single protocol or stack of protocols
- Logic v.s. performance

#### **Tools for Disambiguating RFCs**

Jane Yen Ramesh Govindan University of Southern California yeny@usc.edu ramesh@usc.edu

#### Abstract

For decade, dafting laterest protocols has taken signifcant amounts of human supervision due to the fundamental ambiguty of natural language. Given such ambiguty, it is also not surprised that protocol implementations have long exhibited bags. This pain and overhead can be significantly revelor with the human bags of the simulation processing OUSP specified extensors in BFCs, and to generate protocol implementations automatically when the analyough is clarified. However this systems in far from generat or delongwale. To further reduce the overhead and errors due to ambiguous sentences, and to improve the generality of this system, much work remains to be dues. In this paper, we consider what uexaming the natural-language challenges in the BFC process. COM Conversion.

CCS Concepts

 $\bullet \ Networks \rightarrow Formal \ specifications.$ 

#### Keywords

natural language, protocol specifications ACM Reference Format:

Jane Yen, Ramesh Govindan, and Barath Raghavan. 2021. Tools for Disambiguating RFCs. In Applied Networking Research Workshop (ANRW '21), July 34-30, 2021. Virtual Event. ISAA. ACM. New York, NY, USA, 7 pages. https://doi.org/10.1145/3472345.3472314

#### 1 Introduction

It has long been the case that protocol behaviors are discussed and debated in the networking community for years before they are codified. The laternet Architeture Board (IAB) [9], Internet Engineering Task Force (IETF), and Internet Research Task Force (IETF) [21] canable groups of independent networking professionals to draft Request for

Ramesh Govindan Barath Raghavan University of Southern California ramesh@usc.edu barathra@usc.edu

Commoti (BFC) to explain their networking lease in English. Their aim is to publish their idea globaly and, typically, to colify those ideas as a networking standard. Draft RFCs are carefully evaluated and reviewed by members of working groups and editors, these individuals manually consider editorial and technical perspectives. Toogh the process by which RFCs are developed is effective, and can identify significant vages, increte, or parall discriptions and get authors to their drafts, the overall overhead of this human supervision is significant.

RFC editors follow a number of guidelines to review RFC drafts and provide feedback to the authors: authors typically aim to satisfy the guidelines as guickly as possible so that the back-and-forth review process can be minimized. The guide lines mostly address what sections are required to be present. what order the sections should be in, what minimal information should be stated clearly, and so forth. While these guidelines help drafts across time and space to maintain a uniform style, these guidelines are mostly writing advice and do not effectively help authors or reviewers to identify technical issues. To tackle technical concerns in any RFC draft requires participants' professional background knowledge RFC authors have many ways to describe technical details. including natural language, pseudocode, formal specifica tions, real code, state machine diagrams, and more. While pseudocode, formal specifications, and real code can provide precise execution steps, natural language is still preferable due to its flexibility and readability. However, natural lan guage can also be fuzzy, which sometimes leads to a fundamental technical problem. For example, "the type code changed to 0" is a verbatim sentence in Internet Control Message Protocol (ICMP) RFC [23]. The noun phrase "type code" can confuse a reader as to whether it refers to as a type named "code" or a code named "type" or a specific protocol header field named "type code". If any reader interprets the noun phrase incorrectly, a packet may be generated incorrectly and get dropped by a receiver.

- Paragraph or sentence-based analysis
- Mis-matched/mis-captured behaviors
- Standalone or multiple RFCs
- Single protocol or stack of protocols
- Logic v.s. performance

#### **Tools for Disambiguating RFCs**

Iane Yen University of Southern California venv@usc.edu

#### Abstract

For decades, drafting Internet protocols has taken significant amounts of human supervision due to the fundamental ambiguity of natural language. Given such ambiguity, it is also not surprising that protocol implementations have long exhibited bugs. This pain and overhead can be significantly reduced with the help of natural language processing (NLP). We recently applied NLP to identify ambiguous or underspecified sentences in RFCs, and to generate protocol implementations automatically when the ambiguity is clarified. However this system is far from general or deployable. To further reduce the overhead and errors due to ambiguous sentences, and to improve the generality of this system, much work remains to be done. In this paper, we consider what it would take to produce a fully-general and useful system for easing the natural-language challenges in the RFC process.

#### **CCS** Concepts

Networks → Formal specifications.

#### Keywords

natural language, protocol specifications ACM Reference Format:

Jane Yen, Ramesh Govindan, and Barath Raghavan. 2021. Tools for Disambiguating RFCs. In Applied Networking Research Workshop (ANRW '21), July 24-30, 2021, Virtual Event, USA. ACM, New York, NY, USA, 7 pages. https://doi.org/10.1145/3472305.3472314

#### 1 Introduction

It has long been the case that protocol behaviors are discussed and debated in the networking community for years before they are codified. The Internet Architecture Board (IAB) [9], Internet Engineering Task Force (IETF), and Internet Research Task Force (IRTF) [12] enable groups of independent networking professionals to draft Request for

![](_page_43_Picture_17.jpeg)

This work is licensed under a Creative Commons Attribution in Work is incensed under a Creative Commons Admin inCommercial-ShareAlike International 4.0 License. ANRW '21, July 24-30, 2021, Virtual Event, USA @ 2021 Copyright held by the owner/author(s). ACM 15BN 978-1-4503-8618-8021/07 https://doi.org/10.1145/3472305.3472314

Ramesh Govindan Barath Raghavan University of Southern California University of Southern California ramesh@usc.edu barathra@usc.edu

Comments (RFCs) to explain their networking ideas in English. Their aim is to publish their ideas globally and, typically to codify those ideas as a networking standard. Draft RFCs are carefully evaluated and reviewed by members of working groups and editors; these individuals manually consider editorial and technical perspectives. Though the process by which RFCs are developed is effective, and can identify significant vague, incorrect, or partial descriptions and get authors to their drafts, the overall overhead of this human supervision is significant

RFC editors follow a number of guidelines to review RFC drafts and provide feedback to the authors: authors typically aim to satisfy the guidelines as guickly as possible so that the back-and-forth review process can be minimized. The guide lines mostly address what sections are required to be present. what order the sections should be in, what minimal information should be stated clearly, and so forth. While these guidelines help drafts across time and space to maintain a uniform style, these guidelines are mostly writing advice and do not effectively help authors or reviewers to identify technical issues. To tackle technical concerns in any RFC draft requires participants' professional background knowledge RFC authors have many ways to describe technical details. including natural language, pseudocode, formal specifica tions, real code, state machine diagrams, and more. While pseudocode, formal specifications, and real code can provide precise execution steps, natural language is still preferable due to its flexibility and readability. However, natural lan guage can also be fuzzy, which sometimes leads to a fundamental technical problem. For example, "the type code changed to 0" is a verbatim sentence in Internet Control Message Protocol (ICMP) RFC [23]. The noun phrase "type code" can confuse a reader as to whether it refers to as a type named "code" or a code named "type" or a specific protocol header field named "type code". If any reader interprets the noun phrase incorrectly, a packet may be generated incorrectly and get dropped by a receiver.

- Paragraph or sentence-based analysis
- Mis-matched/mis-captured behaviors
- Standalone or multiple RFCs
- Single protocol or stack of protocols
- Logic v.s. performance

#### **Tools for Disambiguating RFCs** ramesh@usc.edu

Iane Yen Ramesh Govindan University of Southern California venv@usc.edu

#### Abstract

For decades, drafting Internet protocols has taken significant amounts of human supervision due to the fundamental ambiguity of natural language. Given such ambiguity, it is also not surprising that protocol implementations have long exhibited bugs. This pain and overhead can be significantly reduced with the help of natural language processing (NLP). We recently applied NLP to identify ambiguous or underspecified sentences in RFCs, and to generate protocol implementations automatically when the ambiguity is clarified. However this system is far from general or deployable. To further reduce the overhead and errors due to ambiguous sentences, and to improve the generality of this system, much work remains to be done. In this paper, we consider what it would take to produce a fully-general and useful system for easing the natural-language challenges in the RFC process.

**CCS** Concepts

Networks → Formal specifications.

#### Keywords

natural language, protocol specifications ACM Reference Format:

Jane Yen, Ramesh Govindan, and Barath Raghavan. 2021. Tools for Disambiguating RFCs. In Applied Networking Research Workshop (ANRW '21), July 24-30, 2021, Virtual Event, USA. ACM, New York, NY, USA, 7 pages. https://doi.org/10.1145/3472305.3472314

#### 1 Introduction

It has long been the case that protocol behaviors are discussed and debated in the networking community for years before they are codified. The Internet Architecture Board (IAB) [9], Internet Engineering Task Force (IETF), and Internet Research Task Force (IRTF) [12] enable groups of independent networking professionals to draft Request for

This work is licensed under a Creative Commons Attribution in Work is incensed under a Creative Commons Admin inCommercial-ShareAlike International 4.0 License. ANRW '21, July 24-30, 2021, Virtual Event, USA @ 2021 Copyright held by the owner/author(s). ACM 15BN 978-1-4503-8618-8021/07 https://doi.org/10.1145/3472305.3472314

Barath Raghavan University of Southern California University of Southern California barathra@usc.edu

Comments (RFCs) to explain their networking ideas in English. Their aim is to publish their ideas globally and, typically to codify those ideas as a networking standard. Draft RFCs are carefully evaluated and reviewed by members of working groups and editors; these individuals manually consider editorial and technical perspectives. Though the process by which RFCs are developed is effective, and can identify significant vague, incorrect, or partial descriptions and get authors to their drafts, the overall overhead of this human supervision is significant

RFC editors follow a number of guidelines to review RFC drafts and provide feedback to the authors: authors typically aim to satisfy the guidelines as guickly as possible so that the back-and-forth review process can be minimized. The guide lines mostly address what sections are required to be present. what order the sections should be in, what minimal information should be stated clearly, and so forth. While these guidelines help drafts across time and space to maintain a uniform style, these guidelines are mostly writing advice and do not effectively help authors or reviewers to identify technical issues. To tackle technical concerns in any RFC draft requires participants' professional background knowledge RFC authors have many ways to describe technical details. including natural language, pseudocode, formal specifica tions, real code, state machine diagrams, and more. While pseudocode, formal specifications, and real code can provide precise execution steps, natural language is still preferable due to its flexibility and readability. However, natural lan guage can also be fuzzy, which sometimes leads to a fundamental technical problem. For example, "the type code changed to 0" is a verbatim sentence in Internet Control Message Protocol (ICMP) RFC [23]. The noun phrase "type code" can confuse a reader as to whether it refers to as a type named "code" or a code named "type" or a specific protocol header field named "type code". If any reader interprets the noun phrase incorrectly, a packet may be generated incorrectly and get dropped by a receiver.

- Paragraph or sentence-based analysis
- Mis-matched/mis-captured behaviors
- Standalone or multiple RFCs
- Single protocol or stack of protocols
- Logic v.s. performance

#### **Tools for Disambiguating RFCs** ramesh@usc.edu

Iane Yen Ramesh Govindan University of Southern California venv@usc.edu

#### Abstract

For decades, drafting Internet protocols has taken significant amounts of human supervision due to the fundamental ambiguity of natural language. Given such ambiguity, it is also not surprising that protocol implementations have long exhibited bugs. This pain and overhead can be significantly reduced with the help of natural language processing (NLP). We recently applied NLP to identify ambiguous or underspecified sentences in RFCs, and to generate protocol implementations automatically when the ambiguity is clarified. However this system is far from general or deployable. To further reduce the overhead and errors due to ambiguous sentences, and to improve the generality of this system, much work remains to be done. In this paper, we consider what it would take to produce a fully-general and useful system for easing the natural-language challenges in the RFC process.

**CCS** Concepts

Networks → Formal specifications.

Keywords

natural language, protocol specifications ACM Reference Format:

Jane Yen, Ramesh Govindan, and Barath Raghavan. 2021. Tools for Disambiguating RFCs. In Applied Networking Research Workshop (ANRW '21), July 24-30, 2021, Virtual Event, USA. ACM, New York, NY, USA, 7 pages. https://doi.org/10.1145/3472305.3472314

#### 1 Introduction

It has long been the case that protocol behaviors are discussed and debated in the networking community for years before they are codified. The Internet Architecture Board (IAB) [9], Internet Engineering Task Force (IETF), and Internet Research Task Force (IRTF) [12] enable groups of independent networking professionals to draft Request for

This work is licensed under a Creative Commons Attribution in Work is incensed under a Creative Commons Admin inCommercial-ShareAlike International 4.0 License. ANRW '21, July 24-30, 2021, Virtual Event, USA @ 2021 Copyright held by the owner/author(s). ACM 15BN 978-1-4503-8618-8021/07 https://doi.org/10.1145/3472305.3472314

Barath Raghavan University of Southern California University of Southern California barathra@usc.edu

Comments (RFCs) to explain their networking ideas in English. Their aim is to publish their ideas globally and, typically to codify those ideas as a networking standard. Draft RFCs are carefully evaluated and reviewed by members of working groups and editors; these individuals manually consider editorial and technical perspectives. Though the process by which RFCs are developed is effective, and can identify significant vague, incorrect, or partial descriptions and get authors to their drafts, the overall overhead of this human supervision is significant

RFC editors follow a number of guidelines to review RFC drafts and provide feedback to the authors: authors typically aim to satisfy the guidelines as guickly as possible so that the back-and-forth review process can be minimized. The guide lines mostly address what sections are required to be present, what order the sections should be in, what minimal information should be stated clearly, and so forth. While these guidelines help drafts across time and space to maintain a uniform style, these guidelines are mostly writing advice and do not effectively help authors or reviewers to identify technical issues. To tackle technical concerns in any RFC draft requires participants' professional background knowledge RFC authors have many ways to describe technical details. including natural language, pseudocode, formal specifica tions, real code, state machine diagrams, and more. While pseudocode, formal specifications, and real code can provide precise execution steps, natural language is still preferable due to its flexibility and readability. However, natural lan guage can also be fuzzy, which sometimes leads to a fundamental technical problem. For example, "the type code changed to 0" is a verbatim sentence in Internet Control Message Protocol (ICMP) RFC [23]. The noun phrase "type code" can confuse a reader as to whether it refers to as a type named "code" or a code named "type" or a specific protocol header field named "type code". If any reader interprets the noun phrase incorrectly, a packet may be generated incorrectly and get dropped by a receiver.

- Paragraph or sentence-based analysis
- Mis-matched/mis-captured behaviors
- Standalone or multiple RFCs
- Single protocol or stack of protocols
- Logic v.s. performance

#### **Tools for Disambiguating RFCs**

Iane Yen University of Southern California venv@usc.edu

#### Abstract

For decades, drafting Internet protocols has taken significant amounts of human supervision due to the fundamental ambiguity of natural language. Given such ambiguity, it is also not surprising that protocol implementations have long exhibited bugs. This pain and overhead can be significantly reduced with the help of natural language processing (NLP). We recently applied NLP to identify ambiguous or underspecified sentences in RFCs, and to generate protocol implementations automatically when the ambiguity is clarified. However this system is far from general or deployable. To further reduce the overhead and errors due to ambiguous sentences, and to improve the generality of this system, much work remains to be done. In this paper, we consider what it would take to produce a fully-general and useful system for easing the natural-language challenges in the RFC process.

#### **CCS** Concepts

Networks → Formal specifications.

#### Keywords

natural language, protocol specifications ACM Reference Format:

Jane Yen, Ramesh Govindan, and Barath Raghavan. 2021. Tools for Disambiguating RFCs. In Applied Networking Research Workshop (ANRW '21), July 24-30, 2021, Virtual Event, USA. ACM, New York, NY, USA, 7 pages. https://doi.org/10.1145/3472305.3472314

#### 1 Introduction

It has long been the case that protocol behaviors are discussed and debated in the networking community for years before they are codified. The Internet Architecture Board (IAB) [9], Internet Engineering Task Force (IETF), and Internet Research Task Force (IRTF) [12] enable groups of independent networking professionals to draft Request for

![](_page_46_Picture_17.jpeg)

This work is licensed under a Creative Commons Attribution in Work is incensed under a Creative Commons Admin inCommercial-ShareAlike International 4.0 License. ANRW '21, July 24-30, 2021, Virtual Event, USA @ 2021 Copyright held by the owner/author(s). ACM 15BN 978-1-4503-8618-8021/07 https://doi.org/10.1145/3472305.3472314

Ramesh Govindan Barath Raghavan University of Southern California University of Southern California ramesh@usc.edu barathra@usc.edu

> Comments (RFCs) to explain their networking ideas in English. Their aim is to publish their ideas globally and, typically to codify those ideas as a networking standard. Draft RFCs are carefully evaluated and reviewed by members of working groups and editors; these individuals manually consider editorial and technical perspectives. Though the process by which RFCs are developed is effective, and can identify significant vague, incorrect, or partial descriptions and get authors to their drafts, the overall overhead of this human supervision is significant

RFC editors follow a number of guidelines to review RFC drafts and provide feedback to the authors: authors typically aim to satisfy the guidelines as guickly as possible so that the back-and-forth review process can be minimized. The guide lines mostly address what sections are required to be present. what order the sections should be in, what minimal information should be stated clearly, and so forth. While these guidelines help drafts across time and space to maintain a uniform style, these guidelines are mostly writing advice and do not effectively help authors or reviewers to identify technical issues. To tackle technical concerns in any RFC draft requires participants' professional background knowledge RFC authors have many ways to describe technical details. including natural language, pseudocode, formal specifica tions, real code, state machine diagrams, and more. While pseudocode, formal specifications, and real code can provide precise execution steps, natural language is still preferable due to its flexibility and readability. However, natural lan guage can also be fuzzy, which sometimes leads to a fundamental technical problem. For example, "the type code changed to 0" is a verbatim sentence in Internet Control Message Protocol (ICMP) RFC [23]. The noun phrase "type code" can confuse a reader as to whether it refers to as a type named "code" or a code named "type" or a specific protocol header field named "type code". If any reader interprets the noun phrase incorrectly, a packet may be generated incorrectly and get dropped by a receiver.

# **Current Work**

### **Reduce Human Effort**

![](_page_48_Figure_1.jpeg)

• Can we avoid writing an ambiguous sentence in the first place?

• What kind of protocols are we going to support?

### **Solution Directions**

• Can we avoid writing an ambiguous sentence in the first place?

A user interface guides spec author to produce only essential information

• What kind of protocols are we going to support?

### **Solution Directions**

• Can we avoid writing an ambiguous sentence in the first place?

A user interface guides spec author to write unambiguous sentences

What kind of protocols are we going to support?

Stateful protocols (It requires to keep internal states to decide operations)

Name	Description
Packet Format	Packet anatomy (i.e., field structure)
Field Descriptions	Packet header field descriptions
Constraints	Constraints on field values
Protocol Behaviors	Reactions to external/internal events
System Architecture	Protocol implementation components
+ State Management	Session information and/or status
Comm. Patterns	Message sequences (e.g., handshakes)

Table 1: Protocol specification components. SAGE supports those marked with  $\blacklozenge$  (fully) and + (partially).

### Our vision

![](_page_52_Picture_1.jpeg)

User interface

Essential protocol elements

### Our vision

![](_page_53_Picture_1.jpeg)

User interface

Essential protocol elements

(Search) [txt html pdf with_errata bibtex] From: <u>draft_ietf-bfd-base-l1</u> Updated by: <u>7412, 7880, 8562</u>	[Tracker] [NG] (Email) (Diffi) Proposed Standard IPE declarations Errata exist	(Diff2)	(Nite
Internet Engineering Task Force (IETF) Request for Comments: 5880 Category: Standards Track ISSN: 2070-1721	D. Katz D. Ward Juniper Networks June 2010		

#### Bidirectional Forwarding Detection (BFD)

#### Abstract

This document describes a protocol intended to detect faults in the bidirectional path between two forwarding engines, including interfaces, data link[s], and to the eatent possible the forwarding engines themselves, with potentially very low latency. It operates independently of media, data protocols, and routing protocols.

Status of This Memo

This is an Internet Standards Track document.

This document is a product of the Internet Engineering Task Force (IETF). It represents the commensus of the IETF community. It has received public review and has been approved for publication by the Internet Engineering Steering Group (IESG). Further information on Internet Standards is available in <u>Section 1 of RF 371</u>.

Information about the current status of this document, any errats, and how to provide feedback on it may be obtained at http://www.rfc-editor.org/infc/ffc3H80.

Copyright Notice

Copyright (c) 2010 IETF Trust and the persons identified as the document authors. All rights reserved.

This decounts is subject to QC\_JD and the IGT Trust's Legal (http://trustes.iei.org/scilence.ici) as fifted on the date of publication of this document. Decame review them documents carefully, as they describe your rights and restrictions with respect to hims Signification Signature as a second to be documents and the second signature as a second on the second signature to hims Signification Signature as described in the Second signature described in the Signification Signature as a second so the described in the Signification Signature as a second so the described in the Signification Signature as a second so the second signature as a second so the Signature as a second so the described in the Signification Signature as a second so the Signature as a described in the Signification Signature as a second so the Signature as a described in the Signification Signature as a second so the Signature as a described in the Signification Signature as a second so the Signature as a described in the Signification Signature as a second so the Signature as a described in the Signification Signature as a second so the Signature as a described in the Signature as a second so the Signature as a second so the Signature as a described in the Signature as a second so the Signature as a second so the Signature as a described in the Signature as a second so the Signature as a described in the Signature as a second so the Signature as a second

#### **English RFC**

	int time value) /
-	and cype_value, {
	ar -uaca = (char -) (hur + r),
- ''.	8 for echo message;
	U for echo reply message
na	r->type = type_value;
	Set code to 0
hd	r->code = 0;
- //	If code equals 0, an identifier may be zero to help match echos and replie
if	(hdr->code == 0) {
	hdr->identifier = 0;
}	
- 11	If code equals 0, a sequence number may be zero to help match echos and
- 11	replies
if	(hdr->code == 0) {
	hdr->sequence_number = 0;
}	
11	For computing the checksum , the checksum field should be zero
hd	r->checksum = 0;
11	For computing the checksum, if the total length is odd, the received data
11	is padded with one octet of zeros
if	(isodd(length)) {
	pad(édata, sizeof(*data), 0, 1);
	length += 1;
÷	
11	The checksum is the 16-bit one's complement of the one's complement sum of
11	the ICMP message starting with the ICMP Type
- 14	

Executable code

### Our vision

Microsoft   MakeCode Arcade		*	8	٩
	<ul> <li>Setter</li> <li>Press:</li> <li>Store</li> <li>Test size:</li> <li>Test</li></ul>			

User interface

Essential protocol elements

Last received packet

Output packet

**Program states** 

Timer

Transfer (112) [Incl. [

#### Bidirectional Forwarding Detection (BFD)

Abstract

This document describes a protocol intended to detect faults in the bidirectional path between two forwarding engines, including interfaces, data link(s), and to the extent possible the forwarding engines themselves, with potentially very low latency. It operates independently of media, data protocols, and routing protocols.

Status of This Memo

This is an Internet Standards Track document.

This document is a product of the Internet Esgineering Task Force (IETF). It represents the commensus of the IETF community. It has received public review and has been approved for publication by the Internet Engineering Steering Group (IESG). Further information on Internet Standards is available in <u>Section 1 of SRT 574</u>.

Information about the current status of this document, any errata, and how to provide feedback on it may be obtained at http://www.rfc-editor.org/info/rfc5880.

Copyright Notice

Copyright (c) 2010 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is analysis to 200. J and the ITTT Trust's Legal (http://intake.int.com/intake.int/intake.int/intake.int/intake.int/intake.int/intake.int/intake.int/intake.int/intake.int/intake.intak

#### **English RFC**

- F	<pre>void fill_icmp_echo_sender(Echo_or_Echo_Reply_Message_hdr *hdr, uint16_t length,</pre>
I	int type_value) {
I	<pre>char *data = (char *) (hdr + 1);</pre>
I	// 8 for echo message;
I	// 0 for echo reply message
	hdr->type = type_value;
I	// Set code to 0
I	hdr->code = 0;
I	// If code equals 0, an identifier may be zero to help match echos and replies
I	if (hdr->code == 0) {
I	hdr->identifier = 0;
	)
$\mathbf{N}$	// If code equals 0, a sequence number may be zero to help match echos and
- X	// replies
- 1	if (hdr->code == 0) {
I	hdr->sequence_number = 0;
I	}
I	// For computing the checksum , the checksum field should be zero
I	hdr->checksum = 0;
I	// For computing the checksum, if the total length is odd, the received data
I	// is padded with one octet of zeros
I	if (isodd(length)) {
I	<pre>pad(&amp;data, sizeof(*data), 0, 1);</pre>
I	<pre>length += 1;</pre>
I	}
I	// The checksum is the 16-bit one's complement of the one's complement sum of
I	// the ICMP message starting with the ICMP Type
I	hdr->checksun = ul6bit_ones_complement(

Executable code

Q & A

![](_page_55_Picture_1.jpeg)

Yu-Chuan Yen YENY@USC.EDU

SAGE

![](_page_55_Picture_3.jpeg)

https://github.com/USC-NSL/sage.git