

CONTENT MODERATION ON END-TO-END ENCRYPTED SYSTEMS: A LEGAL ANALYSIS

Charles Duan* & James Grimmelmann**

Online messaging platforms like Signal and Google's Messages increasingly use end-to-end encryption (E2EE), in which messages are encrypted on the sender's device and decrypted on the recipient's, so that no one else—not even the platform itself—can read them. Although E2EE protects privacy and advances human rights, the law enforcement community and others have criticized its growing use. In their view, E2EE prevents platforms and government authorities from responding to abuses and criminal activity, including child exploitation, malware, scams, and disinformation. At times, they have argued that E2EE is inherently incompatible with effective content moderation.

Computer science researchers have responded to this challenge with a suite of technologies that enable content moderation on E2EE platforms. These technologies—message franking, forward tracing, homomorphic encryption, and automated client-side scanning—preserve some of the essential privacy guarantees of E2EE while enabling the targets of abuse to detect and report it. These technical advances, however, raise legal questions. If E2EE messages are supposed to be private from a messaging platform, and the platform participates in detecting whether those messages are abusive, is that an “interception” of an “electronic communication” prohibited under the Wiretap Act?

This Article analyzes these new E2EE content moderation technologies in light of six major federal communication statutes: the Wiretap Act, the Stored Communications Act, the Pen Register Act, the Computer Fraud and Abuse Act, the Communications Assistance for Law Enforcement Act, and the PROTECT Our Children Act. While generally we find that these content moderation technologies would pass muster under these statutes, the answers are not as clear-cut as one might hope. The advanced cryptographic techniques that these new content moderation strategies employ raise multiple unsettled questions of law under the communication privacy regimes considered. This legal uncertainty arises not because of the ambiguous ethical nature of the technologies themselves, but because the decades-old statutes failed to accommodate, or indeed contemplate, the innovations in cryptography that enable content moderation to coexist with encryption. To the extent that platforms are limited in their ability to moderate end-to-end encrypted content, then, those limits may arise not from the technology but from the law.

TABLE OF CONTENTS

| | |
|--|----|
| INTRODUCTION | 3 |
| I. BACKGROUND | 9 |
| A. Encryption Technologies | 9 |
| B. History of the “Crypto Wars” | 14 |
| II. COMMUNICATION PRIVACY LAWS | 17 |
| A. The Wiretap Act | 17 |
| B. The Stored Communications Act | 19 |
| C. Pen Registers and Trap-and-Trace Devices | 23 |
| D. The Computer Fraud and Abuse Act | 25 |
| E. Communications Assistance for Law Enforcement Act | 26 |
| F. PROTECT Our Children Act of 2008 | 29 |
| III. E2EE CONTENT MODERATION PROPOSALS | 30 |
| A. Message Franking | 31 |
| 1. <i>Technical Overview</i> | 32 |
| 2. <i>Wiretap Act Analysis</i> | 35 |
| 3. <i>SCA Analysis</i> | 41 |
| 4. <i>PRA Analysis</i> | 43 |
| 5. <i>CFAA Analysis</i> | 45 |
| 6. <i>CALEA Analysis</i> | 47 |
| 7. <i>POCA Analysis</i> | 49 |
| B. Forward Tracing | 50 |
| 1. <i>Technical Overview</i> | 51 |
| 2. <i>Wiretap Act Analysis</i> | 54 |
| 3. <i>SCA Analysis</i> | 55 |
| 4. <i>CALEA Analysis</i> | 57 |
| 5. <i>PRA Analysis</i> | 58 |
| 6. <i>CFAA Analysis</i> | 58 |
| 7. <i>POCA Analysis</i> | 59 |

* Assistant Professor of Law, American University Washington College of Law.

** Tessler Family Professor of Digital and Information Law, Cornell Law School and Cornell Tech. We presented earlier versions of this Article at the 2023 Privacy Law Scholars Conference and the 2023 Computer Science and the Law Scholarship Roundtable at the University of Pennsylvania Carey School of Law. Our thanks to the organizers and participants, and to Aislinn Black, Chris Conley, Thomas Haley, Jason Hartline, Gus Hurwitz, Thomas Kadri, Paul Ohm, Gabe Rudin, Sarah Scheffler, David Thaw, Danny Weitzner, Shane Witnov, Alessia Zornetta, students in the Spring 2023 Law of Software seminar at Cornell Tech, and the editors of the *Georgetown Law Technology Review*. This Article may be freely reused under the terms of the Creative Commons Attribution 4.0 International License, <https://creativecommons.org/licenses/by/4.0>.

| | | |
|-----|---|----|
| C. | Server-Side Automated Content Scanning | 59 |
| 1. | <i>Technical Background</i> | 59 |
| 2. | <i>Wiretap Act Analysis</i> | 61 |
| 3. | <i>SCA Analysis</i> | 63 |
| 4. | <i>PR/TT Analysis</i> | 63 |
| 5. | <i>CALEA Analysis</i> | 64 |
| 6. | <i>CFAA Analysis</i> | 64 |
| 7. | <i>POCA Analysis</i> | 66 |
| D. | Client-Side Automated Content Scanning | 67 |
| 1. | <i>Technical Overview</i> | 68 |
| 2. | <i>Wiretap Act Analysis</i> | 71 |
| 3. | <i>SCA Analysis</i> | 72 |
| 4. | <i>PRA Analysis</i> | 74 |
| 5. | <i>CALEA Analysis</i> | 75 |
| 6. | <i>CFAA Analysis</i> | 76 |
| 7. | <i>POCA Analysis</i> | 78 |
| IV. | DISCUSSION | 80 |
| A. | Statutory Ambiguities and Proposed Amendments | 80 |
| 1. | <i>Information and Content</i> | 81 |
| 2. | <i>Consent and Authorization</i> | 83 |
| 3. | <i>Permitted Business Activities</i> | 84 |
| 4. | <i>Computer Devices</i> | 85 |
| 5. | <i>Making the Statutes More Modular</i> | 86 |
| 6. | <i>CALEA Encryption Exception</i> | 87 |
| B. | Insights into the Technologies | 88 |
| C. | What Is End-to-End Encryption? | 89 |
| | CONCLUSION..... | 92 |

INTRODUCTION

Encryption has costs. Most obviously, there are technical costs. Encrypting and decrypting messages takes time and computing power, and creating secure encrypted systems takes immense engineering effort. Most controversially, there are policy costs. Law enforcement groups object when encryption works as intended, because it makes it harder for authorities to read suspects' communications. And most subtly, there are safety costs. Encrypting messages makes it harder to protect users.

On modern communications platforms, content moderation plays a central role in keeping users safe from spam, harassment, and abuse.¹ To moderate content, a platform must know what that content

¹ See Kate Klonick, *The New Governors: The People, Rules, and Processes Governing Online Speech*, 131 HARV. L. REV. 1598, 1637–39 (2018).

is. But when content is encrypted so that not even the platform itself can read it—when it is protected with *end-to-end encryption*, or E2EE for short—standard techniques of content moderation become impossible. Humans cannot read the messages to see if they contain threats of violence; computers cannot scan them to see if they contain child sexual abuse material (CSAM).

In short, it has appeared that encryption and moderation are incompatible. Heightening the irony, encryption itself is also a safety technology, because privacy is a form of safety.² Encryption advances interests such as data security, privacy, free speech, free association, and other constitutional and human rights.³ Besides undermining these benefits, weakening encryption to enable third-party content access (either for content moderation or law enforcement) creates potentially serious security vulnerabilities.⁴ So

² A. Michael Froomkin & Zak Colangelo, *Privacy as Safety*, 95 WASH. L. REV. 141 (2020).

³ See, e.g., Orin S. Kerr, *The Fourth Amendment in Cyberspace: Can Encryption Create a “Reasonable Expectation of Privacy?”*, 33 CONN. L. REV. 503, 503–04 (2001); A. Michael Froomkin, *Metaphor Is the Key: Cryptography, the Clipper Chip, and the Constitution*, 143 U. PA. L. REV. 709, 810–43 (1995); Christopher Soghoian, *Caught in the Cloud: Privacy, Encryption, and Government Back Doors in the Web 2.0 Era*, 8 J. ON TELECOMMS. & HIGH TECH. L. 359, 375–76, 398–99 (2010); Jan H. Samoriski, John L. Huffman & Denise M. Trauth, *Encryption and the First Amendment*, 2 COMM’N L. & POL’Y 417 (1997); Geoffrey Gordon, Note, *Breaking the Code: What Encryption Means for the First Amendment and Human Rights*, 32 COLUM. HUM. RTS. L. REV. 477 (2001); Gwynne B. Barrett, Note, *Law of Diminishing Privacy Rights: Encryption Escrow and the Dilution of Associational Freedoms in Cyberspace*, 15 N.Y. L. SCH. J. HUM. RTS. 115 (1998); Sean J. Edgett, *Double-Clicking on Fourth Amendment Protection: Encryption Creates a Reasonable Expectation of Privacy*, 30 PEPP. L. REV. 339 (2003).

⁴ See, e.g., Hal Abelson, Ross Anderson, Steven Michael Bellovin, Josh Benaloh, Matt Blaze, Whitfield Diffie, John Gilmore, Peter G. Neumann, Ronald L. Rivest, Jeffrey I. Schiller & Bruce Schneier, *The Risks of Key Recovery, Key Escrow, and Trusted Third-Party Encryption*, 2 WORLD WIDE WEB J. 241, 250–53 (1997); Peter Swire & Kenesa Ahmad, *Encryption and Globalization*, 13 COLUM. SCI. & TECH. L. REV. 416, 432–33 (2012); cf. RICHARD A. CLARKE, MICHAEL J. MORELL, GEOFFREY R. STONE, CASS R. SUNSTEIN & PETER SWIRE, LIBERTY AND SECURITY IN A CHANGING WORLD: REPORT AND RECOMMENDATIONS OF THE PRESIDENT’S REVIEW GROUP ON INTELLIGENCE AND COMMUNICATIONS TECHNOLOGIES 216–19 (2013), https://obamawhitehouse.archives.gov/sites/default/files/docs/2013-12-12_rg_final_report.pdf [<https://perma.cc/ER2Z-KV9G>] (recommending that the United States “not in any way subvert, undermine, weaken, or make vulnerable generally available commercial software” for

the fact that E2EE enhances one form of safety (privacy) while undermining another (protection from abuse) seems like a tragic but inevitable tradeoff.

In the last few years, however, computer science researchers have shown that encryption and moderation are not so incompatible after all. Research teams around the world have developed ways to support content moderation and abuse prevention that do not require letting a communications platform view the unencrypted contents of messages. Indeed, this is such a fruitful area of research that there are already works in the technical literature that taxonomize and systematize research on so-called “content moderation on end-to-end encrypted systems.”⁵ Some of the interest driving this research comes from platforms themselves,⁶ while other research is motivated by the search for technical tools to help address policy problems.⁷

One of these techniques, known in the literature as “message franking,” allows the recipient of an abusive message to report it to a moderator *with receipts*.⁸ By virtue of clever construction of the

encryption).

⁵ E.g., Sarah Scheffler & Jonathan Mayer, *SoK: Content Moderation for End-to-End Encryption*, PROC. ON PRIV. ENHANCING TECHS. 403 (2023), <https://petsymposium.org/popets/2023/popets-2023-0060.php> [<https://perma.cc/222V-5KUZ>]; SENY KAMARA, MALLORY KNODEL, EMMA LLANSÓ, GREG NOJEIM, LUCY QIN, DHANARAJ THAKUR & CAITLIN VOGUS, CTR. FOR DEMOCRACY & TECH., OUTSIDE LOOKING IN: APPROACHES TO CONTENT MODERATION IN END-TO-END ENCRYPTED SYSTEMS (2021), <https://cdt.org/insights/outside-looking-in-approaches-to-content-moderation-in-end-to-end-encrypted-systems/> [<https://perma.cc/8YN9-NE2A>].

⁶ See, e.g., Will Cathcart, WhatsApp, *Encryption Has Never Been More Essential—Or Threatened*, WIRED (Apr. 5, 2021, 9:00 AM), <https://www.wired.com/story/opinion-encryption-has-never-been-more-essential-or-threatened/> [<https://perma.cc/FF6Z-QTGB>] (“[B]y employing sophisticated techniques to analyze metadata, user reports, and other unencrypted information, [WhatsApp] ban[s] millions of dangerous accounts every year.”).

⁷ *Everything in Moderation?*, HORIZON DIGIT. ECON. RSCH. (June 6, 2022), <https://www.horizon.ac.uk/everything-in-moderation/> [<https://perma.cc/6W38-86ZF>] (“E2E encryption presents challenges in dealing with misinformation, disinformation, potentially harmful or illegal content, and striking a balance with freedom of speech.”); Ariadna Matamoros-Fernández, *Encryption Poses Distinct New Problems: The Case of WhatsApp*, 9 INTERNET POL’Y REV. 7, 9 (2020) (“The pervasiveness of encrypted platforms in mediating everyday life in some parts of the world is a reminder that viable content moderation measures without breaking encryption are needed.”).

⁸ See *infra* Part III.A.

receipts, message franking guarantees that recipients can prove that the message they are reporting is authentic, that recipients cannot submit false reports, and that no one besides the recipient of a message can report a message—or even learn anything about it. An extension of message franking, called “forward tracing,” allows the platform to trace a reported message back to its original sender, even if it has been forwarded repeatedly—but again without compromising the privacy of messages that are not reported.⁹

Another broad class of new techniques involves automated scanning of messages to detect problematic content, such as unsolicited photographs of genitalia.¹⁰ Again, it seems like a paradox: how can a platform scan encrypted content? Here, the answer lies in a class of algorithms called “homomorphic encryption.”¹¹ The idea behind homomorphic encryption is that an untrusted party can perform computations on content without knowing what the content is. Imagine a blindfolded chef wearing thick mittens, who follows instructions to take things out of a box, chop them up, put them in the oven for an hour at 350 degrees, and then put them back in the box. This chef can roast vegetables for you but cannot learn whether you were roasting potatoes or parsnips. Similarly, homomorphic encryption allows a platform to run a bad-content detector on a message and report the result to the recipient, without the platform itself learning anything about the content or the result.¹²

Despite the recent explosion in research on moderation in the presence of encryption, legal scholarship has given almost no attention to these new technologies.¹³ As of October 8, 2023, a search of HeinOnline for “message franking” returns no results. Only a small handful of position papers have discussed these technologies’

⁹ See *infra* Part III.B.

¹⁰ See *infra* Part III.C.

¹¹ See *infra* text accompanying notes 314–19.

¹² See DAVID WONG, REAL-WORLD CRYPTOGRAPHY § 15.2, fig.15.6 (2021).

¹³ An exception is client-side scanning, because that technology has received substantial attention in the press. See, e.g., Timothy Gernand, *Scanning iPhones to Save Children: Apple’s on-Device Hashing Algorithm Should Survive a Fourth Amendment Challenge*, 127 DICK. L. REV. 307, 319–20 (2022); Nicholas A. Weigel, *Apple’s “Communication Safety” Feature for Child Users: Implications for Law Enforcement’s Ability to Compel iMessage Decryption*, 25 STAN. TECH. L. REV. 210, 216–17 (2022). Nevertheless, the client-side scanning protocols in the computer science literature go far beyond the relatively simple image-flagging proposals that the legal scholarship has considered. See *infra* Part III.D.

implications for law and policy.¹⁴

This omission is unfortunate because the assumption that encryption and moderation are impossible has been baked into the long-running debates about whether and how to regulate the use of E2EE. Some commentators, for example, have criticized platforms' decisions to deploy E2EE by decrying the consequences for reduced content moderation.¹⁵

¹⁴ See, e.g., Ian Levy & Crispin Robinson, *Thoughts on Child Safety on Commodity Platforms* (July 21, 2022) (unpublished manuscript), <https://arxiv.org/abs/2207.09506> [<https://perma.cc/S63B-RV4W>]; Hal Abelson, Ross Anderson, Steven M. Bellovin, Josh Benaloh, Matt Blaze, Jon Callas, Whitfield Duffie, Susan Landau, Peter G. Neumann, Ronald L. Rivest, Jeffrey I Schiller, Bruce Schneier, Vanessa Teague & Carmela Troncoso, *Bugs in Our Pockets: The Risks of Client-Side Scanning* (Oct. 15, 2021) (unpublished manuscript), <https://arxiv.org/abs/2110.07450> [<https://perma.cc/J55L-GVT4>]; Jonathan Mayer, *Content Moderation for End-to-End Encrypted Messaging* (Oct. 6, 2019) (unpublished discussion paper), https://www.cs.princeton.edu/~jrmayer/papers/Content_Moderation_for_End-to-End_Encrypted_Messaging.pdf [<https://perma.cc/4Y4K-LMQW>]; LINDSEY ANDERSEN, DUNSTAN ALLISON-HOPE & MICHAELA LEE, *BUS. FOR SOC. RESP., HUMAN RIGHTS IMPACT ASSESSMENT: META'S EXPANSION OF END-TO-END ENCRYPTION* (2022), <https://www.bsr.org/reports/bsr-meta-human-rights-impact-assessment-e2ee-report.pdf> [<https://perma.cc/S47B-PRDP>]; see also Gurshabad Grover, Tanaya Rajwade & Divyank Katira, *The Ministry and the Trace: Subverting End-to-End Encryption*, 14 NUJS L. REV. 223 (2021) (considering legality of forward tracing in view of India's constitutional right to privacy).

¹⁵ See Siva Vaidhyanathan, *Be Careful Taking Sides in Mark Zuckerberg vs. William Barr*, SLATE (Oct. 4, 2019, 3:41 PM), <https://slate.com/technology/2019/10/facebook-encryption-mark-zuckerberg-william-barr.html> [<https://perma.cc/8GCL-XQCW>] (suggesting that Facebook has “motivations to install encryption” so that “Facebook can’t be held responsible for failing to keep its system free of calls for violence, harassment, or hate speech”); Natasha Lomas, *UK Tells Messaging Apps Not to Use E2E Encryption for Kids’ Accounts*, TECHCRUNCH (June 30, 2021, 8:43 AM), <https://techcrunch.com/2021/06/30/uk-tells-messaging-apps-not-to-use-e2e-encryption-for-kids-accounts/> [<https://perma.cc/Q66T-UQ5C>] (noting U.K. government guidance to platforms that “[e]nd-to-end encryption makes it more difficult for you to identify illegal and harmful content occurring on private channels”). But see Mike Masnick, *The DOJ Is Conflating the Content Moderation Debate with the Encryption Debate: Don’t Let Them*, TECHDIRT (Oct. 8, 2019, 9:23 AM), <https://www.techdirt.com/2019/10/08/doj-is-conflating-content-moderation-debate-with-encryption-debate-dont-let-them/>

This Article aims to bridge this gap between the computer science and legal literatures. It makes three contributions.

First, we introduce and explain the current universe of content moderation technologies for E2EE platforms. We focus on technologies directed to encrypted communications, rather than techniques that work independently of encryption, such as metadata analysis or platform affordance modifications.¹⁶ Some of the technologies we discuss work around the encryption, either by moderating content automatically on users' devices¹⁷ or by manipulating the encrypted content in ways that have predictable effects on the content despite not revealing it.¹⁸ Others take advantage of users to flag content for moderators' review. The thrust of these technologies is to provide the platform with certainty about who to take action against in response to a valid user report of abuse. The key technical challenge is to provide that certainty without compromising the confidentiality of unreported messages.¹⁹

Second, we provide a careful legal analysis of these moderation technologies *vis à vis* the federal communication privacy laws. Application of these laws to the content moderation technologies in question, we find, is not always straightforward. While in many cases the statutory definitions map well onto the communicative aspects of the technologies, there are often ambiguities in the statutes and caselaw, leaving it unclear how a court would rule. In the worst case, a court might hold that these types of content moderation are actually *illegal*, perversely putting platforms and users back to the stark and unpleasant choice between encryption and moderation. We wish we could say that U.S. communication privacy laws clearly rule out such an outcome, but unfortunately, they do not.

Third, we use our legal analysis as a basis to critique both the statutes and the technologies. Unresolved statutory ambiguities highlight potential areas for reform, such as the need to bring communication privacy laws enacted almost half a century ago in line with more recent developments in cryptographic research. But those ambiguities also highlight potential areas of ethical concern with the content moderation technologies themselves. After all, the communication privacy laws are intended to reflect, however accurately, intuitive and societal norms of privacy. Discrepancies

[<https://perma.cc/L8FS-QRHU>].

¹⁶ See Matamoros-Fernández, *supra* note 7, at 7–8. For example, to limit the rapid spread of disinformation, WhatsApp introduced limits on the number of times messages could be forwarded. See *id.* at 8.

¹⁷ See *infra* Part III.D.

¹⁸ See *infra* Part III.C.

¹⁹ See *infra* Parts III.A, III.B.

between the statutes and the technologies are a starting point for a larger conversation on how those technologies impact privacy interests, an especially significant conversation given the immense privacy benefits of end-to-end encryption.

Part I of the Article gives a technical and historical background on encryption technologies and how they came to pose content moderation challenges.²⁰ Part II is a legal overview. It describes the six most relevant statutes: (1) the Wiretap Act, (2) the Stored Communications Act (SCA), (3) the Pen Register Act (PRA), (4) the Computer Fraud and Abuse Act (CFAA), (5) the Communications Assistance for Law Enforcement Act (CALEA), and (6) the PROTECT Our Children Act of 2008 (POCA).²¹ Part III, the heart of the Article, reviews four principal techniques for content moderation in the presence of E2EE: (1) message franking, (2) forward tracing, and automated scanning with homomorphic encryption either on (3) the platform's servers, or (4) on the user's devices. For each technique, it gives a technical overview, and then analyzes how that technique fares under the various communication privacy statutes.²² Part IV then steps back to extract broader lessons for legal scholars and policy makers.²³

I. BACKGROUND

A. Encryption Technologies

This Part provides a brief overview of the fundamentals of modern cryptography. The reader who is already familiar with the concepts of public-key encryption and hash functions should feel free to skip ahead to the next Part.²⁴

"Encryption" is a process that keeps information confidential by rendering it unintelligible to outsiders. For a simple example, consider ROT-13 encryption, in which every letter in a text is replaced with the letter that is thirteen away in the alphabet:

²⁰ See *infra* Part I.

²¹ See *infra* Part II.

²² See *infra* Part III.

²³ See *infra* Part IV.

²⁴ For more thorough overviews of central concepts in cryptography, see generally (in ascending order of detail) JAMES GRIMMELMANN, INTERNET LAW: CASES AND PROBLEMS 40–45 (12th ed. 2022); NAT'L ACAD. OF SCIS., ENG'G & MED., CRYPTOGRAPHY AND THE INTELLIGENCE COMMUNITY: THE FUTURE OF ENCRYPTION 16–34 (2022); MIKE ROSULEK, THE JOY OF CRYPTOGRAPHY (2021), <https://joyofcryptography.com./pdf/book.pdf> [<https://perma.cc/EMM3-FURT>].

ABCDEFGHIJKLMNOPQRSTUVWXYZ

⇓

NOPQRSTUVWXYZABCDEFGHIJKLM

"WE THE PEOPLE" $\xrightarrow{\text{Encrypt}_{\text{ROT-13}}}$ "JR GUR CRBCYR"

The unencrypted data ("WE THE PEOPLE") is called the "plaintext," and the encrypted data ("JR GUR CRBCYR") is called the "ciphertext."²⁵ To an unauthorized party who does not know how the data has been encrypted, the ciphertext should appear to be random. But an authorized party who knows that it has been encrypted using ROT-13 can recover the plaintext by undoing the letter-by-letter replacement:

"JR GUR CRBCYR" $\xrightarrow{\text{Decrypt}_{\text{ROT-13}}}$ "WE THE PEOPLE"

ROT-13 is not a very good encryption algorithm because its security collapses as soon as an eavesdropper (typically called an "adversary" in the cryptography literature) learns what algorithm is being used. A better approach, and the one which is universally used today, is to use an algorithm that combines the plaintext with an additional piece of information, called an encryption "key," to produce the ciphertext.²⁶ That way, even if the algorithm is publicly known, data encrypted with that algorithm will be indiscernible to an adversary so long as the key is kept appropriately secret.²⁷

ROT-13 is weak because it has no separate key, but a closely related encryption algorithm is stronger because it does. In a "Caesar cipher," each letter in the alphabet is replaced with the letter that is k letters ahead of it in the alphabet (wrapping around at the end, so that A follows Z).²⁸ The number k is the key for a Caesar cipher; it can have any value from 1 to 26. The substitution for a Caesar cipher with a key of 1 is:

²⁵ See ROSULEK, *supra* note 24, at 10.

²⁶ See *id.* at 11.

²⁷ See Orin S. Kerr & Bruce Schneier, *Encryption Workarounds*, 106 GEO. L.J. 989, 993 (2018) (describing Kerchoff's Principle, specifying that "[a]n encryption algorithm should be secure if everything is known about it except the key").

²⁸ See Dennis Luciano & Gordon Prichett, *Cryptology: From Caesar Ciphers to Public-Key Cryptosystems*, 18 COLL. MATH. J. 2, 3–4 (1987).

ABCDEFGHIJKLMNOPQRSTUVWXYZ

⇓

BCDEFGHIJKLMNOPQRSTUVWXYZA

The Caesar cipher with a key of 2 is:

ABCDEFGHIJKLMNOPQRSTUVWXYZ

⇓

CDEFGHIJKLMNOPQRSTUVWXYZAB

and so on. ROT-13 is just a Caesar cipher with a hardwired key of 13.

Caesar ciphers are a form of “symmetric-key” encryption because the key used to decrypt an encrypted ciphertext is the same as or easily derivable from the encryption key.²⁹ To encrypt a message using a Caesar cipher, shift every letter in the plaintext forward by k letters. To decrypt the message, shift every letter in the ciphertext backward by k letters. Symmetric-key encryption can be simple, fast, and convenient, but it also has some significant disadvantages. One of them is the problem of key distribution. Every pair of people who wishes to communicate must coordinate in advance to agree on a secret key to use—and make sure that they do not accidentally reveal the key to an adversary in the process.³⁰

To overcome this, “asymmetric” or “public-key” encryption ciphers use a pair of keys, called the “public key” and the “private key.”³¹ The sender encrypts the message using the public key; the receiver decrypts the message using the private key.

Plaintext $\xrightarrow{\text{Encrypt}_{\text{Public Key}}}$ Ciphertext $\xrightarrow{\text{Decrypt}_{\text{Private Key}}}$ Plaintext

As the names suggest, in many widely used encryption schemes, the public key is truly public, and the private key is truly private. A person who wants to receive messages will generate a key pair and

²⁹ ROSULEK, *supra* note 24, at 260.

³⁰ *See id.* at 12 (noting the difficulty of “key distribution”).

³¹ *See id.* at 260.

widely distribute the public key so that anyone can use it to encrypt messages to them. But the person will keep the private key to themselves, so that they are the only person who can decrypt those messages.³²

This asymmetry is what makes E2EE possible. Suppose that Alice and Bob want to exchange a message on a platform they do not trust. Alice encrypts the message using Bob's public key, and then hands the message off to the platform to deliver to Bob. When Bob receives Alice's message, he can decrypt it using his private key. But because Bob has not shared his private key with anyone, not even the platform can decrypt Alice's message.

By contrast, in a non-E2EE messaging system, the platform has access to the plaintext. This does not mean there is no encryption involved. The communication from Alice to the platform might be encrypted, and so might the communication from the platform to Bob. Furthermore, the platform might encrypt the message when it is "at rest" in storage to keep it safe against hackers. The crucial point, however, is that any encryption that takes place involves keys the platform has access to. It can use those keys whenever it wants, so Alice and Bob must trust the platform not to misuse that power (or be compelled to misuse it). But when Alice encrypts the message to Bob herself using Bob's public key, Alice and Bob only need to trust each other, not the platform.

Public-key encryption is surprisingly versatile. In addition to protecting the privacy of a message, it can be used to establish that the message is authentic. Abstractly, the way that these "digital signatures" work is that if Alice wants to prove her authorship of a message, she encrypts it with her *private* key.³³ Bob can then use Alice's *public* key (which he knows because Alice has shared it with the world) to decrypt the message. Now he knows, to a high degree of certainty, that only Alice could have sent the message because only Alice had access to the private key used to encrypt it.³⁴

The final concept in this brief tour of encryption is the

³² Crucially, it is not feasible to derive the private key from the public key. This is what distinguishes asymmetric encryption from symmetric encryption like a Caesar cipher, where the decryption key is just 26 minus the encryption key.

³³ "Authorship" here is used just to mean that the signer wishes to be identified as associated with the message. It does not mean that they necessarily authored the message in the sense that the author of a novel creates its text and has a copyright in it.

³⁴ More advanced schemes combine standard public-key encryption and digital signatures to ensure that the message is both confidential and signed.

“cryptographic hash,”³⁵ an algorithm that takes input data of arbitrary size, such as a long message, and generates a much smaller “hash value.” A well-designed hash algorithm has several useful features that collectively mean that a hash value uniquely identifies the data it came from without giving away the data itself, similar to how a fingerprint mark uniquely identifies a person but does not reveal the person’s eye color or height. (Hashes are sometimes called “digital fingerprints” for this reason.) Specifically, a good cryptographic hash satisfies the following properties:

- Uniqueness: A given piece of data should predictably produce exactly one hash value.
- Preimage resistance: The hash value should reveal no information about the data, such that one cannot reconstruct the data from the hash.
- Collision resistance: Two different pieces of data should rarely produce the same hash value.

A hash is like a digital signature in that it provides an authenticity guarantee, except that in this case, only someone who had access to the complete text of the original message could have generated the corresponding hash.

Cryptographic hashes have a variety of uses. For one thing, they simplify the digital signature process: a message sender can encrypt just a hash of a message rather than an entire message, and others can still use the signed hash to verify the message’s authorship since the hash uniquely identifies the message. Hashes can also be used to make commitments without publicly revealing the details of what one is committing to. For example, a basketball fan could publish a hash of their predictions for the March Madness NCAA bracket at the start of the tournament, and then reveal their actual bracket after the tournament is over. And finally, hashes can be used in place of content that, for whatever reason, one does not wish to store. For example, major online platforms typically scan uploaded images to see whether users are uploading known examples of CSAM. But, for obvious reasons, a platform does not want to maintain a database filled with images of children being sexually exploited. By storing only the hashes of those images, the platform can still compare uploaded images to known examples of CSAM but without the legal

³⁵ The term “one-way hash” means the same thing. A “hash function” is the mathematical process that transforms input data into a hash, and when the hash function is applied to a message, the result is sometimes called a “message digest” because the hash function has “digested” the message.

and operational nightmares of storing actual CSAM. As we will see, hashes are particularly useful for content moderation in E2EE systems because they allow platforms and users to make verifiable claims about content without revealing the content itself.

B. History of the “Crypto Wars”

Since the introduction of modern asymmetric encryption algorithms in the 1970s, law enforcement and government intelligence agencies have argued that widespread private use of encryption would hamper criminal investigations and national security efforts.³⁶ But for many cryptographers and privacy activists, government surveillance was the principal threat that made widespread use of cryptography a moral necessity. The policy debates between these two camps for and against laws restricting the use of encryption were informally dubbed the “Crypto Wars.”

In the 1990s, these debates came to a head in the United States. The federal government for a time deemed strong encryption a “munition” subject to export restrictions; lawmakers proposed “key-escrow” technologies, in which government agencies would hold special decryption keys that would enable law enforcement to decrypt communications upon receipt of a court order.³⁷ These skirmishes ended with defeats for the government. The export controls were relaxed to allow academic cryptographers and open-source programmers to post their work online without fear of prosecution, and the leading proposed key-escrow scheme, the Clipper Chip, collapsed in ignominy when it was shown to be insecure.

But although the 1990s left encryption legal and widely used, it was not omnipresent. In particular, most communications platforms were not *end-to-end* encrypted. A message from Alice to Bob would be encrypted in transit from Alice to the platform, encrypted in

³⁶ See generally CRAIG JARVIS, *CRYPTO WARS: THE FIGHT FOR PRIVACY IN THE DIGITAL AGE* 111–52 (2021); STEVEN LEVY, *CRYPTO: HOW THE CODE REBELS BEAT THE GOVERNMENT—SAVING PRIVACY IN THE DIGITAL AGE* (2001). For a broader history of cryptography, see generally DAVID KAHN, *THE CODEBREAKERS: THE COMPREHENSIVE HISTORY OF SECRET COMMUNICATION FROM ANCIENT TIMES TO THE INTERNET* (1996).

³⁷ See DANIELLE KEHL, ANDI WILSON & KEVIN BANKSTON, *DOOMED TO REPEAT HISTORY? LESSONS FROM THE CRYPTO WARS OF THE 1990S* 5–12 (2015), <http://newamerica.org/cybersecurity-initiative/policy-papers/doomed-to-repeat-history-lessons-from-the-crypto-wars-of-the-1990s/> [<https://perma.cc/7GYB-VM77>]; see generally LEVY, *supra* note 36.

storage on the platform, and encrypted in transit from the platform to Bob—but the platform held the decryption keys that would allow it to decrypt the message received from Alice, and to retrieve and decrypt the messages it stored. This meant that law enforcement could still effectively obtain unencrypted communications between platform users—either by serving the platform with legal process or by infiltrating the platform’s systems to copy out the data and decryption keys.

Following Edward Snowden’s revelations of widespread national security surveillance on Internet communications, technology companies increasingly rearchitected their communications platforms to thwart this surveillance. If the platform itself was a potential vulnerability, the natural solution was end-to-end encryption. Now the message from Alice to Bob would be encrypted using a key pair controlled by Bob, so that the platform would have no greater ability to read the message than any random stranger would. This development led to a new wave of complaints in the mid-2010s from law enforcement about the danger of malfeasants “going dark” and the need for “back doors” to give government access to encrypted content.³⁸ Content moderation on E2EE systems has become a point of contention in this larger debate over strong encryption. There are at least three interlocking problems.

First, and most seriously from law enforcement’s point of view, E2EE makes it harder to detect and investigate the transmission of illegal content, such as CSAM and terrorist plots. This is not only a surveillance and security problem, but also a content moderation problem. Identifying and removing content that has predictably harmful effects for third parties is a classic goal of content moderation.

Second, E2EE makes it harder for platforms to defend users from spam, abuse, and harassment. When the platform has access to all communications, it can flag specific messages that are likely to be unwanted and identify suspicious patterns of mass coordinated messaging. These capabilities appear to disappear when all messages

³⁸ See KEHL, WILSON & BANKSTON, *supra* note 37, at 1; NAT’L ACAD. OF SCIS., ENG’G & MED., *DECRYPTING THE ENCRYPTION DEBATE: A FRAMEWORK FOR DECISION MAKERS* 7–9 (2018); Harold Abelson, Ross Anderson, Steven M. Bellovin, Josh Benaloh, Matt Blaze, Whitfield Diffie, John Gilmore, Matthew Green, Susan Landau, Peter G. Neumann, Ronald L. Rivest, Jeffrey I. Schiller, Bruce Schneier, Michael A. Specter & Daniel J. Weitzner, *Keys Under Doormats: Mandating Insecurity by Requiring Government Access to All Data and Communications*, 1 J. CYBERSECURITY 69 (2015); WHITFIELD DIFFIE & SUSAN LANDAU, *PRIVACY ON THE LINE: THE POLITICS OF WIRETAPPING AND ENCRYPTION* (2d ed. 2010).

are end-to-end encrypted.

Third, content moderation for groups typically depends on delegating some of the moderation work to group administrators and other users.³⁹ These users make moderation decisions that are enforced by the platform. But again, when all group communications are indecipherable to the platform, it appears that it cannot effectively intervene to carry out the instructions of group administrators.

The result, as many commentators have noted, is that encrypted group communications can become vectors for abusive harassment and disinformation campaigns of the sort that platforms regularly moderate.⁴⁰

Thus, the collateral harms to content moderation have become a standard argument against strong end-to-end encryption. In 2019, government officials from the United States, the United Kingdom, and Australia sent an open letter to Mark Zuckerberg, Chief Executive Officer of Facebook, calling on the company and other online platforms to “not deliberately design their systems to preclude any form of access to content”—that is, not to implement end-to-end encryption that would “severely erod[e] a company’s ability to detect and respond to illegal content and activity.”⁴¹ Others have similarly argued that when encryption prevents platforms from reading users’ messages, the platforms are unable to identify or respond to the online sexual abuse of children.⁴²

Privacy advocates and computer scientists, in turn, have

³⁹ James Grimmelman, *The Virtues of Moderation*, 17 YALE J.L. & TECH. 42, 69–70, 95 (2015) (discussing delegation in Reddit).

⁴⁰ See Matamoros-Fernández, *supra* note 7, at 7–8 (citing NIC NEWMAN, RICHARD FLETCHER, ANTONIS KALOGEROPOULOS & RASMUS KLEIS NIELSEN, REUTERS INSTITUTE DIGITAL NEWS REPORT 2019, at 9 (2019), https://reutersinstitute.politics.ox.ac.uk/sites/default/files/2019-06/DNR_2019_FINAL_0.pdf [<https://perma.cc/5Y4H-HL7>]); Cristina Tardáguila, Fabrício Benevenuto & Pablo Ortellado, *Fake News Is Poisoning Brazilian Politics. WhatsApp Can Stop It*, N.Y. TIMES (Oct. 17, 2018, 3:00 PM), <https://www.nytimes.com/2018/10/17/opinion/brazil-election-fake-news-whatsapp.html> [<https://perma.cc/D6B3-ZJB5>].

⁴¹ See Letter from Priti Patel, Secretary of State, U.K. Home Off., William P. Barr, Att’y Gen., U.S. Dep’t of Just., Kevin K. McAleenan, Acting Secretary, U.S. Dep’t of Homeland Sec. & Peter Dutton, Minister, Austl. Dep’t of Home Affs., to Mark Zuckerberg, Chief Exec. Officer, Facebook, *Open Letter: Facebook’s “Privacy First” Proposals* 1 (Oct. 4, 2019), <https://www.justice.gov/opa/press-release/file/1207081/download> [<https://perma.cc/PC64-E4TN>].

⁴² See, e.g., *End-to-End Encryption*, NAT’L CTR. FOR MISSING & EXPLOITED CHILD., <http://www.missingkids.org/e2ee.html> [<https://perma.cc/ZEP9-68QU>].

challenged these criticisms.⁴³ They observed that content moderation is a suite of strategies broader than mere reading of messages: it includes user reporting workflows, automated data analysis, and message flagging mechanisms.⁴⁴ And to substantiate their claim that content moderation can be compatible with E2EE, computer science researchers have redoubled their efforts toward developing novel content moderation strategies.⁴⁵

II. COMMUNICATION PRIVACY LAWS

Having introduced the background and history of E2EE technology, this Part now turns to the communication privacy laws that cover platforms offering such encryption. Our focus is on six major federal regimes: the Wiretap Act, the Stored Communications Act, the Pen Register Act, the Computer Fraud and Abuse Act, the Communications Assistance for Law Enforcement Act, and the PROTECT Our Children Act. The key provisions of each are discussed below.

A. The Wiretap Act

Originally enacted in 1968 and as amended by the Electronic Communications Privacy Act of 1986 (“ECPA”) and other statutes,⁴⁶ § 2511 governs the interception of live communications.⁴⁷ Generally,

⁴³ See Mayer, *supra* note 14; KAMARA ET AL., *supra* note 5.

⁴⁴ See KAMARA ET AL., *supra* note 5, at 7–11; *see also* Grimmelmann, *supra* note 39, at 55–79.

⁴⁵ See, e.g., Anunay Kulshrestha & Jonathan Mayer, *Identifying Harmful Media in End-to-End Encrypted Communication: Efficient Private Membership Computation*, 30 PROC. USENIX SEC. SYMP. 893 (2021) (noting government officials’ letter to Zuckerberg as a motivation for developing privacy-preserving perceptual hash technology). To be sure, some of these content moderation technologies predate the 2019 letter. Facebook developed message franking in 2017. See FACEBOOK, INC., MESSENGER SECRET CONVERSATIONS: TECHNICAL WHITEPAPER 11–12 (ver. 2.0 2017), <https://about.fb.com/wp-content/uploads/2016/07/messenger-secret-conversations-technical-whitepaper.pdf>. [<https://perma.cc/PM9U-WF5P>].

⁴⁶ See Electronic Communications Privacy Act of 1986, Pub. L. No. 99-508, §§ 101–11, 100 Stat. 1848, 1848–59.

⁴⁷ See Wiretap Act, 18 U.S.C. § 2511; *see generally* JIM DEMPSEY, ARI SCHWARTZ & ALISSA COOPER, CTR. FOR DEMOCRACY & TECH., AN OVERVIEW OF THE FEDERAL WIRETAP ACT, ELECTRONIC COMMUNICATIONS PRIVACY ACT, AND STATE TWO-PARTY CONSENT LAWS OF RELEVANCE TO THE NEBUAD SYSTEM AND OTHER USES OF INTERNET TRAFFIC CONTENT FROM ISPs FOR BEHAVIORAL ADVERTISING 3–11

the statute creates liability for one who “intentionally intercepts”⁴⁸ the “contents of any wire, oral, or electronic communication”⁴⁹ by means of an “electronic, mechanical, or other device,”⁵⁰ unless the interception falls under an exception in the statute.⁵¹ The statute is typically characterized as requiring five elements:⁵²

1. *intentional*—“inadvertent conduct is no crime; the offender must have done on purpose those things which are outlawed.”⁵³
2. *interception*—the communication must be “captured or redirected,”⁵⁴ perhaps in a separate, contemporaneous transmission.⁵⁵
3. *of the contents*—“information concerning the substance, purport, or meaning,”⁵⁶ and not mere metadata,⁵⁷ must be received.
4. *of an electronic communication*—including messages both as they are in transit and when in transient electronic storage.⁵⁸
5. *using a device*—“The term includes computers, but it is defined so as not to include hearing aids or extension telephones in normal use”⁵⁹

(2008), <https://cdt.org/wp-content/uploads/privacy/20080708ISPTraffic.pdf> [<https://perma.cc/ZLN9-GDSM>].

⁴⁸ 18 U.S.C. § 2511(1)(a).

⁴⁹ *Id.* § 2510(4).

⁵⁰ *Id.* § 2510(5).

⁵¹ *See id.* § 2511(2).

⁵² *See, e.g., In re Pharmatrak, Inc. Priv. Litig.*, 329 F.3d 9, 18 (1st Cir. 2003); CHARLES DOYLE, CONG. RSCH. SERV., R41733, *PRIVACY: AN OVERVIEW OF THE ELECTRONIC COMMUNICATIONS PRIVACY ACT 7* (ver. 9 2012).

⁵³ DOYLE, *supra* note 52, at 7.

⁵⁴ *United States v. Rodriguez*, 968 F.2d 130, 136 (2d Cir. 1992).

⁵⁵ *See Pharmatrak*, 329 F.3d at 21–22; *United States v. Councilman*, 418 F.3d 67, 80 (1st Cir. 2005) (discussing but not resolving meaning of “contemporaneous”); *cf. Konop v. Hawaiian Airlines, Inc.*, 302 F.3d 868, 878–79 (9th Cir. 2002) (holding no interception occurred based on illicit viewing of website content long after content was posted).

⁵⁶ 18 U.S.C. § 2510(8).

⁵⁷ *See, e.g., United States v. N.Y. Tel. Co.*, 434 U.S. 159, 166–67 (1977); *In re Google Inc. Cookie Placement Consumer Priv.*, 806 F.3d 125, 135–39 (3d Cir. 2015).

⁵⁸ *See Councilman*, 418 F.3d at 79.

⁵⁹ DOYLE, *supra* note 52, at 10.

Several exceptions to § 2511 are relevant. Law enforcement may cause the interception of otherwise protected communications, provided that certain procedural requirements are met.⁶⁰ There is no violation if one of the parties to the communication consents to the interception⁶¹ or if the interceptor is a party to the communication.⁶² Additionally, the provider of a communication service may intercept communications “in the ordinary course of its business,”⁶³ as a “necessary incident” to providing the service,⁶⁴ or for “the protection of the rights or property of the provider.”⁶⁵

B. The Stored Communications Act

While the Wiretap Act deals with interception of data in transit, the SCA deals with access to communications and data in electronic storage.⁶⁶ Enacted as part of ECPA in 1986, the SCA pairs a general prohibition on unauthorized access to stored communications⁶⁷ with specific rules directed to service providers.⁶⁸

A key threshold element of the SCA is whether a communication is in “electronic storage.”⁶⁹ This term is defined narrowly, encompassing only “temporary, intermediate storage” or “storage . . . for purposes of backup protection.”⁷⁰ One question is whether this definition encompasses messages that the recipient has retrieved but that still remain in the service provider’s storage. A line of cases, including *Theofel v. Farey-Jones*, has deemed such post-transmission stored messages as “backup protection,”⁷¹ but other

⁶⁰ See Wiretap Act, 18 U.S.C. § 2511(2)(a)(ii).

⁶¹ See *id.* § 2511(2)(c). Analogous state laws sometimes require all parties to the communication to consent to interception. See DEMPSEY ET AL., *supra* note 47, at 11 & n.37.

⁶² See Wiretap Act, 18 U.S.C. § 2511(2)(d).

⁶³ *Id.* § 2510(5)(a)(ii).

⁶⁴ *Id.* § 2511(2)(a)(i).

⁶⁵ *Id.*

⁶⁶ See generally Orin S. Kerr, *A User’s Guide to the Stored Communications Act, and a Legislator’s Guide to Amending It*, 72 GEO. WASH. L. REV. 1208 (2004) [hereinafter Kerr, *User’s Guide*].

⁶⁷ See Stored Communications Act, 18 U.S.C. § 2701.

⁶⁸ See *id.* § 2702.

⁶⁹ See *id.* § 2701(a) (flush text); *id.* § 2702(a)(1).

⁷⁰ See 18 U.S.C. § 2510(17).

⁷¹ *Theofel v. Farey-Jones*, 359 F.3d 1066, 1077 (9th Cir. 2004); see also *Quon v. Arch Wireless Operating Co.*, 529 F.3d 892, 902 (9th Cir. 2008), *rev’d on other grounds sub nom. City of Ontario v. Quon*, 560 U.S. 746 (2010); *O’Grady v. Superior Ct.*, 44 Cal. Rptr. 3d 72, 84 n.9 (Cal. Ct. App. 2006); *Konop v. Hawaiian Airlines, Inc.*, 302 F.3d 868, 879–80 (9th Cir. 2002) (assuming without deciding that posts on a website are in

courts have questioned that view.⁷² Persistent browser cookies, by contrast, likely do not fall within this definition.⁷³

The general unauthorized access prohibition is contained in 18 U.S.C. § 2701. Under that section, it is a violation to:

- *access without authorization or exceed authorized access*⁷⁴
- *a “facility through which an electronic communication service is provided”*—generally an online server but possibly also a user device⁷⁵—
- to obtain, alter, or prevent access to a *wire or electronic*

“electronic storage”).

⁷² There are two lines of argument here. First, some courts posit that the SCA is limited to communications stored during transit, so post-transmission messages are outside the statute’s ambit. *See, e.g.,* *Fraser v. Nationwide Mut. Ins. Co.*, 135 F. Supp. 2d 623, 636 (E.D. Pa. 2001). Others observe that for services such as webmail, a post-transmission message stored on a web server is not backing up any other copy of the message. *See, e.g.,* *Sartori v. Schrodts*, 424 F. Supp. 3d 1121, 1132–33 (N.D. Fla. 2019) (quoting *United States v. Weaver*, 636 F. Supp. 2d 769, 772 (C.D. Ill. 2009)); *cf. Theofel*, 359 F.3d at 1076 (“An ISP that kept permanent copies of temporary messages could not fairly be described as ‘backing up’ those messages.”).

⁷³ *See In re DoubleClick Inc. Priv. Litig.*, 154 F. Supp. 2d 497, 512–13 (S.D.N.Y. 2001); *In re iPhone Application Litig.*, 844 F. Supp. 2d 1040, 1059 (N.D. Cal. 2012). *In re Intuit Privacy Litigation* denied a motion to dismiss a complaint alleging an SCA violation based on persistent cookies, but the court gave only cursory attention to whether the cookies were in “electronic storage,” and the defendant did not appear to have pressed the issue. 138 F. Supp. 2d 1272, 1275–76 (C.D. Cal. 2001). *Chance v. Avenue A, Inc.* found no SCA violation for accessing persistent cookies on unrelated grounds and never reached the question of whether the cookies were in “electronic storage.” *See* 165 F. Supp. 2d 1153, 1161–62 (W.D. Wash. 2001). In view of these cases, Kerr’s suggestion in *User’s Guide* that “several district courts have applied the SCA to regulate the placement of electronic cookies on home computers” is somewhat puzzling. Kerr, *User’s Guide*, *supra* note 66, at 1214.

⁷⁴ 18 U.S.C. § 2701(a)(1)–(2).

⁷⁵ *Id.* § 2701(a)(1). *Compare DoubleClick*, 154 F. Supp. 2d at 509 (suggesting that a user’s personal computer on which a website cookie is stored is a “facility”), and *Chance*, 165 F. Supp. 2d at 1161 (“[I]t is possible to conclude that modern computers, which serve as a conduit for the web server’s communication . . . , are facilities covered under the Act.”), with Kerr, *User’s Guide*, *supra* note 66, at 1215 & n.47 (arguing that home computers are not electronic communication services).

*communication*⁷⁶

- in *electronic storage*, as above.

Section 2701 further provides exceptions for access authorized by the communications service provider,⁷⁷ by the user who sends or receives the communication,⁷⁸ or by law enforcement.⁷⁹

The specific provisions for service providers in 18 U.S.C. § 2702 deal with divulgence of information stored with service providers. That section creates two modes of liability depending on whether the service provider offers an “electronic communication service” (ECS) or a “remote computing service” (RCS).⁸⁰ The vexed distinction between the two is that an ECS provides its users “the ability to send or receive wire or electronic communications,”⁸¹ while an RCS provides “computer storage or processing services.”⁸² Most authorities agree that a service can satisfy both definitions.⁸³ For example, an online email service might be an electronic communication service when it holds onto an email before the recipient reads it, but a remote computing service after the email is read insofar as the recipient uses the email service for long-term storage.⁸⁴ If a service acts as both at the same time, then it must avoid both sets of prohibitions to escape liability.⁸⁵

⁷⁶ 18 U.S.C. § 2701(a) (flush text).

⁷⁷ *Id.* § 2701(c)(1).

⁷⁸ *Id.* § 2701(c)(2).

⁷⁹ *Id.* § 2701(c)(3) (incorporating the SCA’s provisions for law-enforcement access).

⁸⁰ *See id.* § 2702(a)(1)–(2).

⁸¹ Wiretap Act, 18 U.S.C. § 2510(15).

⁸² Stored Communications Act, § 2711(2).

⁸³ *See, e.g.,* Theofel v. Farey-Jones, 359 F.3d 1066, 1076–77 (9th Cir. 2004); Kerr, *User’s Guide*, *supra* note 66, at 1215 (“[M]ost network service providers . . . can act as providers of ECS in some contexts, providers of RCS in other contexts, and as neither in some contexts as well.”).

⁸⁴ *See* Kerr, *User’s Guide*, *supra* note 66, at 1215–16.

⁸⁵ In *Quon*, the Ninth Circuit held that a text messaging provider was an electronic communication provider liable under § 2702(a)(1) for disclosing a police officer’s messages to the city that employed the officer. *See* *Quon v. Arch Wireless Operating Co.*, 529 F.3d 892, 897–98 (9th Cir. 2008). The city argued that, as the subscriber to the text service, it fell within a statutory exception that only applied to remote computing services. *See id.* at 900. In concluding that the exception did not apply, the Ninth Circuit reasoned that the text service was “more appropriately categorized as an ECS than an RCS.” *Id.* at 902. Although this could be taken to mean that the two categories are mutually exclusive, a better reading is that the

A violation of the statute requires:

- For a public ECS:
 - knowingly divulging;
 - the contents of a communication; and
 - in electronic storage,⁸⁶ or
- For a public RCS:
 - knowingly divulging;
 - content of communication on the service;
 - from (or created for) a subscriber or customer of the service;
 - solely for the purpose of the storage or computing services; and
 - if the provider is not authorized to access the contents for other purposes.⁸⁷

Like the Wiretap Act, § 2702 provides several exceptions to the prohibitions on services' divulging communications. No violation occurs for divulging the content of a communication:

- to the addressee or intended recipient of the communication,⁸⁸
- with the consent of the sender or recipient,⁸⁹
- for forwarding the communication to its destination,⁹⁰
- as “necessarily incident to the rendition of the service[;]”⁹¹

court held that liability could arise based on the text messaging service's electronic communication service capacity regardless of whether it was also a remote computing service.

⁸⁶ See 18 U.S.C. § 2701(a)(1).

⁸⁷ See *id.* § 2701(a)(2).

⁸⁸ See *id.* § 2702(b)(1).

⁸⁹ See *id.* § 2702(b)(3).

⁹⁰ See *id.* § 2702(b)(4).

⁹¹ See *id.* § 2702(b)(5).

- for “protection of the rights or property of the provider[;]”⁹²
or
- to the government under appropriate circumstances.⁹³

A different violation occurs when a public electronic communications service or remote computing service divulges non-content customer information to the government.⁹⁴ The statute offers a more limited set of exceptions to this prohibition. Non-content customer information may be disclosed with consent, as an incident of rendering the service, to protect the service provider’s rights or property, or to law enforcement under appropriate circumstances.⁹⁵

C. Pen Registers and Trap-and-Trace Devices

While the Wiretap Act deals with the acquisition of the *contents* of a communication, the PRA deals with the acquisition of *metadata* about the parties to a communication. Also known as the Pen Register Act, the statute, again enacted as part of ECPA in 1986, prohibits the use of devices that record wire or electronic communications metadata without a court order.⁹⁶ For historical reasons related to telephone technology, the PRA uses two different terms to describe the regulated devices: a “pen register” records information about *outgoing* communications sent from a device,⁹⁷ while a “trap-and-trace device” records information about *incoming* communications sent to a device.⁹⁸ When the distinction between the two is immaterial, we will refer to them collectively as “PR/TT devices.”

For both types of devices, the key definitional phrase is that these devices capture “dialing, routing, addressing, and signaling information” (DRAS) about wire or electronic communications.⁹⁹ Despite their telephone-era names, the definitions of these tools

⁹² See *id.* § 2702(b)(6).

⁹³ See *id.* § 2702(b)(2), (6)–(9).

⁹⁴ See *id.* § 2702(a)(3).

⁹⁵ See *id.* § 2701(c).

⁹⁶ Pen Register Act, 18 U.S.C. § 3121(a) (“[N]o person may install or use a pen register or a trap and trace device without first obtaining a court order . . .”).

⁹⁷ *Id.* § 3127(3).

⁹⁸ *Id.* § 3127(4).

⁹⁹ *Id.* § 3127(3) (defining a pen register as “a device or process which records or decodes [DRAS]”); *id.* § 3127(4) (defining a trap-and-trace device as “a device or process which captures the incoming electronic or other impulses which identify . . . [DRAS].”).

encompass digital-era technologies, and courts have held that systems for capturing addresses in emails,¹⁰⁰ internet protocol (IP) addresses,¹⁰¹ and physical location information¹⁰² can qualify as PR/TT devices.

The definitions of PR/TT devices are limited in several ways. Most importantly, a device designed to capture content is outside the scope of the PRA;¹⁰³ such a device is regulated instead by the Wiretap Act.¹⁰⁴ To qualify as a regulated PR/TT device, it must also collect communication metadata sent in the course of the communication, not at a different time or from a third-party data source.¹⁰⁵ A device used for billing purposes is exempted from the definition of pen registers, but not from the definition of trap-and-trace devices.¹⁰⁶ Finally, at least one court has suggested that a communication recipient's collection of metadata from the communication does not constitute operation of a trap-and-trace device.¹⁰⁷

Furthermore, the statute provides several exceptions for cases in which a provider's operation of a PR/TT device is automatically legal and does not require advance court authorization:

- “relating to the operation, maintenance, or testing” of the service;¹⁰⁸
- relating “to the protection of the rights or property of such provider, or to the protection of users of that service from

¹⁰⁰ See *In re Application of the U.S. for an Ord. Authorizing the Installation & Use of a Pen Reg. & a Trap & Trace Device on E-Mail Acct.*, 416 F. Supp. 2d 13, 16 (D.D.C. 2006).

¹⁰¹ See *United States v. Soybel*, 13 F.4th 584, 590–94 (7th Cir. 2021).

¹⁰² See *United States v. Sanchez-Jara*, 889 F.3d 418, 420 (7th Cir. 2018); *United States v. Patrick*, 842 F.3d 540, 543 (7th Cir. 2016).

¹⁰³ See 18 U.S.C. § 3127(3)–(4).

¹⁰⁴ See, e.g., *In re Innovatio IP Ventures, LLC Pat. Litig.*, 886 F. Supp. 2d 888, 895 (N.D. Ill. 2012). Accidental capture of content may be permissible so long as the operator of the device “takes all reasonably available steps to minimize the collection of content information and is prohibited from making use of any content information that may be collected.” *In re Certified Question of L.*, 858 F.3d 591, 598 (FISA Ct. Rev. 2016). *United States v. Fregoso*, 60 F.3d 1314, 1321 (8th Cir. 1995); *Brown v. Waddell*, 50 F.3d 285, 291 (4th Cir. 1995).

¹⁰⁶ Compare 18 U.S.C. § 3127(3), with *id.* § 3127(4).

¹⁰⁷ See *Capitol Recs., Inc. v. Thomas-Rassett*, No. 06-1497, slip op. at 8 (D. Minn. June 11, 2009) (reasoning that “the Internet could not function” if recipients could not collect metadata).

¹⁰⁸ 18 U.S.C. § 3121(b)(1).

abuse of service or unlawful use of service[;]”¹⁰⁹

- “to record the fact” of a communication in order to protect the provider or a user of the service “from fraudulent, unlawful or abusive use of service[;]”¹¹⁰
- with “the consent of the user.”¹¹¹

D. The Computer Fraud and Abuse Act

First enacted in 1984 and subsequently amended many times, the CFAA prohibits unauthorized trespass into computer systems.¹¹² The statute provides several pathways to a violation, the broadest of which¹¹³ requires:

- intentional access to a *protected computer*—including any computer “used in or affecting interstate or foreign commerce or communication,”¹¹⁴ which covers “every computer connected to the Internet.”¹¹⁵ The term “computer” is defined broadly, and an interconnected network of encrypted messaging users may constitute a “computer” under the statute.¹¹⁶
- *without authorization or exceeding authorized access*, discussed below;
- to cause various harms, such as obtaining information from

¹⁰⁹ *Id.*

¹¹⁰ *Id.* § 3121(b)(2).

¹¹¹ *Id.* § 3121(b)(3).

¹¹² See Computer Fraud and Abuse Act, 18 U.S.C. § 1030; Orin S. Kerr, *Vagueness Challenges to the Computer Fraud and Abuse Act*, 94 MINN. L. REV. 1561, 1565–68 (2010) [hereinafter Kerr, *Vagueness Challenges*] (summarizing history of amendments to the CFAA).

¹¹³ See 18 U.S.C. § 1030(a)(2)(C).

¹¹⁴ *Id.* § 1030(e)(2).

¹¹⁵ Kerr, *Vagueness Challenges*, *supra* note 112, at 1568.

¹¹⁶ See 18 U.S.C. § 1030(e)(1); Jonathon W. Penney & Bruce Schneier, *Platforms, Encryption, and the CFAA: The Case of WhatsApp v. NSO Group*, 36 BERKELEY TECH. L.J. 469, 478–79 (2022); see also Jonathan Mayer, *The “Narrow” Interpretation of the Computer Fraud and Abuse Act: A User Guide for Applying United States v. Nosal*, 84 GEO. WASH. L. REV. 1644, 1653–54 (2016) (noting how cloud computing requires reconceptualizing of “the scope of a computer system”).

the computer,¹¹⁷ obtaining value by fraud,¹¹⁸ or damaging the computer in some cases.¹¹⁹

What constitutes “authorization” has been a central question in CFAA jurisprudence over the years.¹²⁰ Most recently, the Supreme Court in *Van Buren v. United States* adopted a “gates-up-or-down” approach to authorization: a person either has or lacks access to a computer resource, and contractual restrictions on how the resource is to be used do not affect CFAA authorization.¹²¹ The Court entertained but did not adopt a requirement that any limit on authorization be “code-based” via technological access measures.¹²² As a result, contractual and other legal notions of consent can define the scope of authorization under the statute.¹²³ Social norms and expectations can inform what qualifies as authorization,¹²⁴ and commentators have suggested that the use of E2EE can constitute a denial of authorization to third parties without the encryption keys—making E2EE a code-based limit on authorization.¹²⁵

E. Communications Assistance for Law Enforcement Act

Although the Wiretap Act permits law enforcement to intercept communications under appropriate circumstances,¹²⁶ the statute does not guarantee that the communications system will be engineered to allow for such interception. Congress recognized this problem in the early 1990s, as telephone providers transitioned from analog to digital systems incompatible with existing wiretap techniques.¹²⁷

¹¹⁷ See 18 U.S.C. § 1030(a)(2)(C).

¹¹⁸ See *id.* § 1030(a)(4).

¹¹⁹ See *id.* § 1030(a)(5). As distinct from other violations of the CFAA, a violation of § 1030(a)(5) requires access “without authorization,” not simply access that exceeds authorization. See *Int’l Airport Ctrs., LLC v. Citrin*, 440 F.3d 418, 420 (7th Cir. 2006).

¹²⁰ See generally PETER G. BERRIS, CONG. RSCH. SERV., R46536, CYBERCRIME AND THE LAW: COMPUTER FRAUD AND ABUSE ACT (CFAA) AND THE 116TH CONGRESS 6–7 (Sept. 21, 2020).

¹²¹ *Van Buren v. United States*, 141 S. Ct. 1648, 1658–59 (2021).

¹²² See *id.* at 1659 n.8.

¹²³ See James Grimmelmänn, *Consenting to Computer Use*, 84 GEO. WASH. L. REV. 1500, 1502–03 (2016) [hereinafter Grimmelmänn, *Consenting*].

¹²⁴ See Orin S. Kerr, *Norms of Computer Trespass*, 116 COLUM. L. REV. 1143, 1146 (2016) [hereinafter Kerr, *Norms*].

¹²⁵ See Penney & Schneier, *supra* note 116, at 494–95.

¹²⁶ See Wiretap Act, 18 U.S.C. § 2511(2)(a)(ii).

¹²⁷ See, e.g., *U.S. Telecom Ass’n v. Fed. Commc’ns Comm’n*, 227 F.3d 450, 454 (D.C. Cir. 2000); KRISTIN FINKLEA, CONG. RSCH. SERV.,

Balancing law enforcement concerns with policy issues raised by privacy advocates, Congress enacted CALEA in 1994.¹²⁸

Under the statute, telecommunications carriers must build their systems to be capable of isolating and enabling government interception of a subscriber's wire and electronic communications as well as call-identifying information.¹²⁹ The threshold requirement for the statute to apply is whether a service is a "telecommunications carrier."¹³⁰ Initially, the statute defines this term broadly as "a person or entity engaged in the transmission or switching of wire or electronic communications as a common carrier for hire."¹³¹ This may seem capacious enough to encompass any online platform,¹³² but the definition is subject to two limitations. First, the statute provides that "telecommunications carrier" includes electronic communication services "to the extent that the Commission finds that such service is a replacement for a substantial portion of the local telephone exchange service."¹³³ According to the D.C. Circuit, this provision allows the Federal Communications Commission (FCC) "to expand the definition of a 'telecommunications carrier' to include new technologies that substantially replace the functions of an old-fashioned telephone network."¹³⁴ Second, CALEA exempts from its scope "persons or entities insofar as they are engaged in providing information services."¹³⁵ The term "information services" broadly encompasses services offering "capability for generating, acquiring,

R44187, ENCRYPTION AND EVOLVING TECHNOLOGY: IMPLICATIONS FOR U.S. LAW ENFORCEMENT INVESTIGATIONS 2 (2016); Justin (Gus) Hurwitz, *Encryption*^{Congress} *mod* (Apple + CALEA), 30 HARV. J.L. & TECH. 355, 373–76 (2017).

¹²⁸ See Communications Assistance for Law Enforcement Act (CALEA) § 103, 47 U.S.C. § 1002; FINKLEA, *supra* note 127, at 2.

¹²⁹ See CALEA § 103(a).

¹³⁰ See *id.*

¹³¹ See *id.* § 101(8)(A).

¹³² The term "common carrier for hire" is not defined, but CALEA incorporates the definitions of the Wiretap Act, which in turn incorporates the definition of "communication common carrier" from the Communications Act of 1934. See CALEA § 101(1); Wiretap Act, 18 U.S.C. § 2510(10). The Communications Act, in turn, has no definition of "communication common carrier" but does define "common carrier" as "any person engaged as a common carrier for hire, in interstate or foreign communication by wire or radio or interstate or foreign radio transmission of energy," subject to several exceptions. Communications Act of 1934, 47 U.S.C. § 153(11).

¹³³ CALEA § 101(8)(B)(ii).

¹³⁴ Am. Council on Educ. v. Fed. Commc'ns Comm'n, 451 F.3d 226, 228 (D.C. Cir. 2006).

¹³⁵ CALEA § 101(8)(C)(i).

storing, transforming, processing, retrieving, utilizing, or making available information via telecommunications,” and includes “electronic messaging services.”¹³⁶ As a result, CALEA is generally understood not to apply to most online services, other than those that have features similar to traditional telephony.¹³⁷

The FCC is also responsible for adopting technical standards for CALEA’s capabilities requirements, to the extent that industry and law enforcement are unable to agree on standards independently.¹³⁸ CALEA is thus unusual among communications laws in that it gives law enforcement a hand in the technological design of systems ordinarily left to private industry.¹³⁹

CALEA recites three major limitations upon its assistance requirement.¹⁴⁰ First, the statute provides that it neither imposes requirements of specific technological designs nor prohibits adoption of any particular technology.¹⁴¹ Second, the statute reiterates that its requirements do not apply to “information services.”¹⁴² Third, CALEA does not require telecommunications carriers to provide for “decrypting, or ensuring the government’s ability to decrypt, any communication . . . unless the encryption was provided by the carrier and the carrier possesses the information necessary to decrypt the communication.”¹⁴³ Since in a standard E2EE system the “information necessary to decrypt the communication” lies solely with the communicants, this exception applies to platforms offering such encryption.¹⁴⁴

¹³⁶ *Id.* § 101(6)(A), (B)(iii); *see id.* § 102(4) (defining “electronic messaging services” as “software-based services that enable the sharing of data, images, sound, writing, or other information among computing devices controlled by the senders or recipients of the messages”). This definition is distinct from the Communications Act definition of “information service,” and under current interpretations, definitively encompasses Internet service providers. *See Am. Council on Educ.*, 451 F.3d at 232.

¹³⁷ *See Am. Council on Educ.*, 451 F.3d at 232 (upholding Federal Communications Commission determination that voice-over-IP services qualified as telecommunication services under CALEA); Hurwitz, *supra* note 127, at 383–84.

¹³⁸ *See* CALEA § 107(b).

¹³⁹ *See* Susan Landau, *National Security on the Line*, 4 J. ON TELECOMMS. & HIGH TECH. L. 409, 412, 417–18 (2006).

¹⁴⁰ A fourth exception, irrelevant to the present paper, relates to services for private networks and interconnection of carriers. *See* CALEA § 103(b)(2)(B).

¹⁴¹ *See id.* § 103(b)(1).

¹⁴² *See id.* § 103(b)(2)(A).

¹⁴³ *Id.* § 103(b)(3).

¹⁴⁴ *See, e.g., In re Ord. Requiring Apple, Inc. to Assist in the Execution*

F. PROTECT Our Children Act of 2008

POCA,¹⁴⁵ as amended by the CyberTipline Modernization Act of 2018,¹⁴⁶ is the primary framework for online platforms' responsibilities with respect to CSAM. Three of its statutory provisions are potentially relevant to platforms performing content moderation on E2EE systems.

First, 18 U.S.C. § 2258A imposes a duty upon online service providers¹⁴⁷ to report CSAM to the CyberTipline operated by the National Center for Missing and Exploited Children (NCMEC).¹⁴⁸ The provider making the report must also “preserve any visual depictions, data, or other digital files that are reasonably accessible and may provide context or additional information about the reported material or person.”¹⁴⁹

Despite these stringent reporting requirements and the serious penalties for failures to report,¹⁵⁰ the duty that § 2258A imposes is fairly limited. Reporting is mandatory only for providers with “actual knowledge” of a violation,¹⁵¹ and there is no duty for a provider to “monitor any user, subscriber, or customer” or to “affirmatively search, screen, or scan” for violations.¹⁵² Nor does the report have to contain any specific information: the contents of the report are left to “the sole discretion of the provider.”¹⁵³ These caveats mean that § 2258A will have limited effect on platforms adopting end-to-end encryption or content moderation technologies thereon because the encryption in most cases will limit the platform's actual knowledge of violative content.

The second statutory provision of relevance is 18 U.S.C. § 2258C, which provides that NCMEC may distribute “elements” of

of a Search Warrant Issued by this Court, 149 F. Supp. 3d 341, 355 n.13 (E.D.N.Y. 2016); Hurwitz, *supra* note 127, at 381–82.

¹⁴⁵ PROTECT Our Children Act of 2008 (POCA), Pub. L. No. 110-401, 122 Stat. 4229.

¹⁴⁶ CyberTipline Modernization Act of 2018, Pub. L. No. 115-395, 132 Stat. 5287.

¹⁴⁷ See 18 U.S.C. § 2258E(6) (defining “provider” as any “electronic communication service provider or remote computing service”).

¹⁴⁸ *Id.* § 2258A(a)(1)(A), (a)(2)(A)–(B).

¹⁴⁹ *Id.* § 2258A(h)(2).

¹⁵⁰ See *id.* § 2258A(e).

¹⁵¹ *Id.* § 2258A(a)(1)(A).

¹⁵² *Id.* § 2258A(f)(1), (3); see *United States v. Stevenson*, 727 F.3d 826, 830 (8th Cir. 2013) (holding that § 2258A “makes clear that an electronic communication service provider is not required to monitor any user or communication”).

¹⁵³ 18 U.S.C. § 2258A(b).

reports, such as “hash values or other unique identifiers,” to online service providers.¹⁵⁴ Use of NCMEC’s hash database is optional¹⁵⁵ and is strictly limited to “the sole and exclusive purpose of permitting [providers] to stop the online sexual exploitation of children.”¹⁵⁶ To the extent that platforms adopt content moderation technologies based on NCMEC’s hash database to identify improper content, the platforms presumably must comply with this limitation, although it is unclear what the scope of the limitation is and what penalties would lie for noncompliance.¹⁵⁷

Third, 18 U.S.C. § 2258B immunizes providers from any “civil or criminal charge . . . arising from the performance of the reporting or preservation responsibilities” in § 2258A or § 2258C.¹⁵⁸ At first glance, this section might seem to be a useful defense against liability under the other communication privacy laws, to the extent that a platform’s content moderation activities relate to reporting CSAM. Indeed, reporting under § 2258A is an explicit exception to the SCA’s prohibition on divulging stored communications.¹⁵⁹ However, courts have construed this immunity narrowly, holding that it covers the act of disclosing information to NCMEC but not the acts of intercepting or searching for illicit content prior to disclosure.¹⁶⁰

III. E2EE CONTENT MODERATION PROPOSALS

With the background on encryption and the statutory regimes in mind, we can now turn to the slate of new technologies for E2EE content moderation. The technologies that we review are called

¹⁵⁴ *Id.* § 2258C(a)(1)–(2); *see* United States v. Coyne, 387 F. Supp. 3d 387, 400 (D. Vt. 2018) (“Section 2258C specifically authorizes the hash value technology used in PhotoDNA.”); NAT’L CTR. FOR MISSING & EXPLOITED CHILD., OUR 2022 IMPACT 14 (Apr. 2023), <https://www.missingkids.org/content/dam/missingkids/pdfs/2022-ncmec-our-impact.pdf> [<https://perma.cc/8UTY-9CPS>] (describing “hash-sharing list” of “6,314,832 hashes”).

¹⁵⁵ *See* 18 U.S.C. § 2258C(c).

¹⁵⁶ *Id.* § 2258C(a)(1), (b).

¹⁵⁷ As of October 8, 2023, no court appears to have interpreted § 2258C.

¹⁵⁸ 18 U.S.C. § 2258B(a). An exception to this immunity is made for intentional or reckless harmful acts. *See id.* § 2258B(b).

¹⁵⁹ *See* Stored Communications Act, 18 U.S.C. § 2702(b)(6); United States v. Rosenow, 50 F.4th 715, 730 (9th Cir. 2022).

¹⁶⁰ *See* United States v. Stevenson, 727 F.3d 826, 830 (8th Cir. 2013) (“Section 2258B(a) . . . is silent regarding whether or how [a provider] should scan its users’ e-mail.”), *followed in* United States v. DiTomaso, 81 F. Supp. 3d 304, 311 (S.D.N.Y. 2015); United States v. Ackerman, 831 F.3d 1292, 1297 (10th Cir. 2016) (interpreting § 2258B(a) not to authorize providers to “review [content] intentionally”).

“message franking,” “forward tracing,” “server-side scanning,” and “client-side scanning.” For each, we provide a brief overview of the technical operations of the technology as described in the computer science literature, and then proceed to evaluate those operations under each of the six communication privacy laws being considered.

A. Message Franking

Introduced as part of Facebook’s secret-conversations service¹⁶¹ and elaborated upon in research,¹⁶² message franking enables moderation of messages that users report to the platform as abusive or otherwise in violation of the platform’s policies.¹⁶³ A challenge with moderation of user-flagged content is that the platform must be able to verify who sent an illicit message before it can take action against them.¹⁶⁴ For non-E2EE messages, the platform can inspect its logs to confirm that the message was sent by the user it appears to be from. But in a standard E2EE system, the platform has no independent way to verify that a reported message was sent by its putative sender—or that it was sent on the platform at all. The objective of message franking is to tie the content of a message to its sender, avoiding potential forgeries and false accusations, so the platform can take appropriate action against abusers with confidence.¹⁶⁵

Digital signatures have long provided sender verifiability, but message franking aims to achieve a second objective called “deniability” or “repudiability” that digital signatures do not achieve.¹⁶⁶ “Off-the-record” conversations are often important for maintaining separations across social contexts: a person communicating with family and friends may not want to leave a permanent, provable record of those conversations that could wind

¹⁶¹ See FACEBOOK, INC., *supra* note 45, at 11–12.

¹⁶² See, e.g., Paul Grubbs, Jiahui Lu & Thomas Ristenpart, *Message Franking via Committing Authenticated Encryption*, 37 PROC. ANN. INT’L CRYPTOLOGY CONF. 66, 67 (2017) (initiating “formal study of message franking”); Yevgeniy Dodis, Paul Grubbs, Thomas Ristenpart & Joanne Woodage, *Fast Message Franking: From Invisible Salamanders to Encryption*, 38 PROC. ANN. INT’L CRYPTOLOGY CONF. 155 (2018).

¹⁶³ For a general description of message franking, see KAMARA ET AL., *supra* note 5, at 17–18.

¹⁶⁴ See Grubbs et al., *supra* note 162, at 75.

¹⁶⁵ See *id.*

¹⁶⁶ See *id.*; Nikita Borisov, Ian Goldberg & Eric Brewer, *Off-the-Record Communication, or, Why Not to Use PGP*, 2004 PROC. ACM WORKSHOP ON PRIV. ELEC. SOC’Y 77, 79 (2004).

up in the hands of coworkers, bosses, or the world.¹⁶⁷ A digitally-signed message creates that permanent record because anyone with access to the message sender's public key can verify the sender's authorship of the message.¹⁶⁸ Message franking protocols, by contrast, aim for deniability, such that no one else can provably tie message contents to their senders. Only the platform can, and it can do so only when it has received an abuse report from the recipient.¹⁶⁹

Message franking offers a simple but powerful mechanism for user-reported content moderation on E2EE systems. However, the technology requires the content moderator or platform to manipulate the message's content over the wire, potentially raising legal questions.

1. Technical Overview

Message franking consists generally of a series of steps performed by a message sender, the messaging platform, and the recipient. To send a message, the sender first generates a cryptographic hash of the message with a randomly generated secret key.¹⁷⁰ This hash, the "franking tag," is sent to the platform along with the message and the random key, the latter two items (but not the franking tag) being encrypted using the recipient's public key.¹⁷¹ Since the random key is encrypted and thus unreadable to the platform, the franking tag at this point is meaningless to the platform.

¹⁶⁷ See generally HELEN NISSENBAUM, *PRIVACY IN CONTEXT: TECHNOLOGY, POLICY, AND THE INTEGRITY OF SOCIAL LIFE* (2009).

¹⁶⁸ See Borisov et al., *supra* note 166, at 79.

¹⁶⁹ See Grubbs et al., *supra* note 162, at 75.

¹⁷⁰ See FACEBOOK, INC., *supra* note 45, at 11 (*NF* and *TF* are the secret key and cryptographic hash, respectively); Grubbs et al., *supra* note 162, at 75 (*Kf* and *C2*). The variables are identified to help track the respective concepts across the different notations used in the cited references.

¹⁷¹ See FACEBOOK, INC., *supra* note 45, at 11; Grubbs et al., *supra* note 162, at 75.

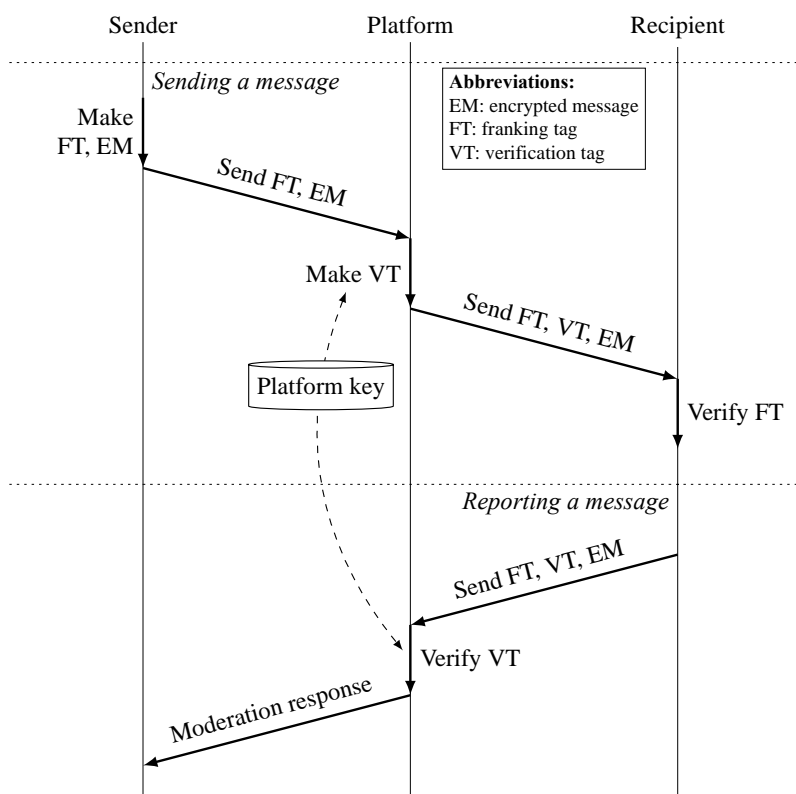


Figure 1: Diagram of communications performed in message franking.

The platform maintains its own platform secret key, and uses that key, the franking tag, and the sender's identification to produce another cryptographic hash, the "verification tag."¹⁷² Because the platform key is secret, third parties cannot forge verification tags that misrepresent the senders of messages. The platform sends the encrypted content to the intended message recipient, who decrypts the content and random key. The recipient verifies that the franking tag was correctly generated; if it was not, the message is discarded as fraudulent.¹⁷³

The verification tag, being hashed using the platform's secret key, is meaningless to third parties, so it cannot be used to breach any sender confidentiality that the platform offers. Nevertheless, if the message recipient flags the message for moderation, then the

¹⁷² See FACEBOOK, INC., *supra* note 45, at 12 (*KF* and *RF* are the platform key and verification tag, respectively); Grubbs et al., *supra* note 162, at 75 (*KFB* and *a*).

¹⁷³ See FACEBOOK, INC., *supra* note 45, at 12 ("If *TF* is not verified then the recipient discards the message without displaying it.").

platform can confirm the sender's identity. To do so, the recipient sends the unencrypted message, random key, sender's identification, and verification tag to the platform.¹⁷⁴ Along with the platform's secret key, the platform now has all the information necessary to recreate the verification tag, thereby proving that the reported message and sender identity were truthful. The platform can now take content moderation actions against the sender of the improper content, confident that the sender in fact sent that improper content.

Improvements to the above message franking system allow for greater anonymity,¹⁷⁵ or enhanced performance and security.¹⁷⁶ In some, though not all, cases, the platform may retain parts of the communication, including the verification tag. First, if the recipient's device is not connected to the platform at the time, then the platform may hold onto the message and tags temporarily in the course of delivery. Second, Facebook's messaging service treats messages with attachments differently: it stores the attachment on the platform's servers to be later downloaded by the recipient, and also stores elements of the message franking communication on its servers.¹⁷⁷

¹⁷⁴ See *id.*; Grubbs et al., *supra* note 162, at 75 (“To report abuse, the recipient sends [the message], *Kf*, and *a* to Facebook.”).

¹⁷⁵ See Nirvan Tyagi, Paul Grubbs, Julia Len, Ian Miers & Thomas Ristenpart, *Asymmetric Message Franking: Content Moderation for Metadata-Private End-to-End Encryption*, 39 PROC. ANN. INT’L CRYPTOLOGY CONF. 222 (2019) (discussing message franking on “sealed-sender” platforms where the platform is unaware of the sender’s identity at the time of message delivery); Long Chen & Qiang Tang, *People Who Live in Glass Houses Should Not Throw Stones: Targeted Opening Message Franking Schemes*, CRYPTOLOGY EPRINT ARCHIVE (Dec. 14, 2018), <https://eprint.iacr.org/2018/994> [<https://perma.cc/9GH8-869W>] (discussing message franking where the recipient need not reveal the entire message contents).

¹⁷⁶ See Rawane Issa, Nicolas Alhaddad & Mayank Varia, *Hecate: Abuse Reporting in Secure Messengers with Sealed Sender*, 31 PROC. USENIX SEC. SYMPOSIUM 2335, 2339 (2022) (adding “forward” and “backward” security to message franking); Hiroki Yamamuro, Keisuke Hara, Masayuki Tezuka, Yusuke Yoshida & Keisuke Tanaka, *Forward Secure Message Franking*, 24 INT’L CONF. ON INFO. SEC. & CRYPTOLOGY 339, 340–41 (2021) (discussing a message franking scheme resilient to compromises of platform keys); Dodis et al., *supra* note 162 (developing more efficient algorithm for message franking).

¹⁷⁷ See FACEBOOK, INC., *supra* note 45, at 9; Dodis et al., *supra* note 162, at 161–62 (identifying a bug in Facebook’s implementation of attachment storage and message franking).

2. Wiretap Act Analysis

Interception of communications implicating the Wiretap Act might occur at three points in the above message franking protocol: (1) when the platform receives the message and franking tag and signs them to generate a verification tag, (2) when the recipient opens the message, and (3) when the platform receives an abuse notification. Of these three points, only the first raises any substantial legal issue. On (2), the recipient opening the message does not violate the statute because they are “a party to the communication.”¹⁷⁸ And on (3), the platform receiving an abuse notification similarly has “prior consent to such interception” from the recipient,¹⁷⁹ and furthermore, the platform merely “listens to or copies the communication that has already been captured,” which courts have held not to constitute interception under the statute.¹⁸⁰

In turn, the question of whether the platform’s reading and signing the franking tag of a message violates the Wiretap Act depends on four issues: (1) whether the franking tag is “content,” (2) whether the platform is the intended recipient of the franking tag, (3) whether the users have consented to the franking, and (4) whether franking is in the platform’s ordinary course of business.

Content. Since the Wiretap Act is limited to interception of content, no violation would occur if the franking tag is not content. Yet that determination is surprisingly difficult. The statute defines “contents” as “any information concerning the substance, purport, or meaning” of a communication.¹⁸¹ The franking tag, being a cryptographic hash of the message, carries no informational value to the platform. Yet, being derived from the content and computed so that virtually no other message would produce the same hash value, the franking tag is inextricably tied to the substance of the message and so, is arguably “information concerning” that message.

While scholars have recognized uncertainty in the meaning of “contents” under the Wiretap Act, the literature has largely focused on the substantive value of metadata and not whether encrypted content is content.¹⁸² The case law is no clearer. The case most on

¹⁷⁸ Wiretap Act, 18 U.S.C. § 2511(2)(d).

¹⁷⁹ *Id.*

¹⁸⁰ *Noel v. Hall*, 568 F.3d 743, 749 (9th Cir. 2009). If the message is an audio message and the abuse notification is sent *before* the recipient listens to the message, then there may be an interception under the statute, as discussed *infra* Part III.D.2.

¹⁸¹ 18 U.S.C. § 2510(8).

¹⁸² *See, e.g., In re Google Inc. Cookie Placement Consumer Priv.*, 806

point applied the Wiretap Act, though only indirectly, to a hash of an illicit file used in a peer-to-peer filesharing network.¹⁸³ The district court found it a “closer question” whether the hash was content, and avoided it by relying on unrelated grounds for its decision.¹⁸⁴ Another district court held that a flag identifying whether an email message was encrypted was content, “[h]owever trivial.”¹⁸⁵ By contrast, a series of decisions about users of pirate devices for decrypting satellite television broadcasts held that no interception of content could occur by mere receipt of encrypted data unless it was decrypted.¹⁸⁶

These competing views suggest that at least some courts, but not all, might be persuaded that an encrypted message qualifies as content under the Wiretap Act.¹⁸⁷ The encrypted message’s

F.3d 125, 137 nn.37–38 (3d Cir. 2015) (quoting 2 WAYNE R. LAFAYE, JEROLD H. ISRAEL, NANCY J. KING & ORIN S. KERR, *CRIMINAL PROCEDURE* § 4.4(d) (3d ed. 2007) (content “depends entirely on the circumstances”); Orin Kerr, *Websurfing and the Wiretap Act*, WASH. POST (June 4, 2015, 1:29 AM), <https://www.washingtonpost.com/news/volokh-conspiracy/wp/2015/06/04/websurfing-and-the-wiretap-act/> [<https://perma.cc/X8U9-BJ5R>] (“[T]he line between contents and metadata is not abstract but contextual with respect to each communication.”); Steven M. Bellovin, Matt Blaze, Susan Landau & Stephanie K. Pell, *It’s Too Complicated: How the Internet Opens Katz, Smith, and Electronic Surveillance Law*, 30 HARV. J.L. & TECH. 1 (2016). The latter article discusses one proposed test for distinguishing content from metadata: If data being transmitted can be encrypted without affecting transport of the data, then the data is content. See Bellovin, Blaze, Landau & Pell, *supra*, at 78–79 (discussing Shane Huang, Distinguishing Content from Metadata: The Provider-Conscious Encryption Test (May 2, 2014) (unpublished student paper)). While the authors dispute that proposed test, it is an indication that at least some scholars might accept encrypted content to be content.

¹⁸³ See *United States v. Sigouin*, 494 F. Supp. 3d 1252, 1263–64 (S.D. Fla. 2019). Specifically, the defendant in the case sought to suppress evidence under the Fourth Amendment, and argued that the standard for unlawful interception under the Wiretap Act should inform the court’s assessment of his reasonable expectation of privacy. See *id.* at 1263.

¹⁸⁴ *Sigouin*, 494 F. Supp. 3d at 1264.

¹⁸⁵ *Optiver Austl. Pty. Ltd. v. Tibra Trading Pty. Ltd.*, No. 12-80242, 2013 WL 256771, at *2 (N.D. Cal. Jan. 23, 2013).

¹⁸⁶ See, e.g., *DirecTV, Inc. v. Barnes*, 302 F. Supp. 2d 774, 779 (W.D. Mich. 2004).

¹⁸⁷ For what it is worth, Congress arguably views encrypted data as content subject to the Wiretap Act, in view of a 2000 amendment to the statute requiring an annual report on the number of wiretap orders “in which encryption was encountered.” Automatic Elimination and Sunset Reports Exemption Act, Pub. L. No. 106-97, § 2(a), 114 Stat. 246, 247

ciphertext provides information about the sender's desire to use encryption, the ciphertext perhaps could be used in comparison against other messages, and the ciphertext certainly would reveal the sender's content if the interceptor later obtained access to the decryption keys. Although some courts might choose to follow the satellite cases in holding unintelligible encrypted messages not to be content, it would be unwise for the designer of a messaging system to assume that all courts would reach that result.

Consent. Even if the platform is not itself a party to the communication, it is possible that the sender or the recipient of the message consented to the platform's interception of the franking tag. Such consent is a defense to Wiretap Act liability,¹⁸⁸ and a message platform's terms of service can suffice as consent to interception.¹⁸⁹ To be sure, courts have been leery of treating broad statements about data use in privacy policies or terms of service as specific consent to interception of communications, and may look unfavorably upon broad, nonspecific terms of service as evidence of users' consent to interception.¹⁹⁰

Another issue is that the platform may not be able to get consent from all users. If a platform accepts messages to or from third-party services not operated by the platform, then the senders or recipients of those messages may not have provided consent to the platform's interception of franking tags or other message information.¹⁹¹ One

(2000) (codified at 18 U.S.C. § 2519, (2)(b)(iv)). The same statute added reporting on PR/TT devices but included nothing on encryption there, again suggesting that Congress viewed encrypted data as content rather than metadata. *See id.* § 3, 114 Stat. at 247–48.

¹⁸⁸ *See* Wiretap Act, 18 U.S.C. § 2511(2)(a).

¹⁸⁹ *See, e.g., In re Yahoo Mail Litig.*, 7 F. Supp. 3d 1016, 1028–31 (N.D. Cal. 2014).

¹⁹⁰ *See, e.g., In re Google Inc. Gmail Litig.*, No. 13-2430, slip op. at 23–27 (N.D. Cal. Sept. 26, 2013); *In re Pharmatrak, Inc. Priv. Litig.*, 329 F.3d 9, 20–21 (1st Cir. 2003) (“Consent ‘should not casually be inferred.’”) (quoting *Griggs-Ryan v. Smith*, 904 F.2d 112, 117–18 (1st Cir. 1990)).

¹⁹¹ *See Gmail*, No. 13-2430, slip op. at 27–28; Bruce E. Boyden, *Can a Computer Intercept Your Email?*, 34 CARDOZO L. REV. 669, 678 (2012) (providing examples where consent may not be obtained). Existing E2EE messaging systems are closed universes that do not interoperate with third-party services, but the European Union's Digital Markets Act will require messaging interoperability that preserves security. *See generally* Julia Len, Esha Ghosh, Paul Grubbs & Paul Rösler, *Interoperability in End-to-End Encrypted Messaging* (2023) (unpublished manuscript), <https://eprint.iacr.org/2023/386.pdf> [<https://perma.cc/SA3K-KQMN>] (discussing technical challenges in implementing E2EE-preserving messaging interoperability).

might argue that the act of transmitting the franking tag is implicit consent to the platform's use of the tag. Because the recipient verifies the correctness of the franking tag before accepting a message for delivery, the sender cannot have a message delivered and read without producing a valid franking tag.¹⁹² The sender's request to have the message delivered, then, arguably entails authorization to generate the verification tag.

Courts have found implicit consent to interception under the Wiretap Act,¹⁹³ but have often been reluctant to do so.¹⁹⁴ As a result, it is not clear whether consent to interception under the Wiretap Act would be found if explicit user agreement is absent. Such situations may arise with more frequency should the technological environment move toward interoperable messaging systems.¹⁹⁵

Intended Recipient. Even if the franking tag is content, the platform would be permitted to intercept it in transit if the platform can show that it was “a party to the communication.”¹⁹⁶

Several cases illustrate the complexity of this intended-recipient exception. In *In re Google Inc. Cookie Placement Consumer Privacy*, a class of website visitors accused Google of unlawful interception when the websites contained code instructing users' web browsers to transmit user information to Google.¹⁹⁷ The Third Circuit concluded there was no such interception because the information was sent directly by a web request from the visitors' browsers to Google.¹⁹⁸

¹⁹² See FACEBOOK, INC., *supra* note 45, at 12.

¹⁹³ See *In re DoubleClick Inc. Priv. Litig.*, 154 F. Supp. 2d 497, 510 (S.D.N.Y. 2001) (inferring consent from “technological and commercial relationships with its affiliated Web sites”).

¹⁹⁴ See *Pharmatrak*, 329 F.3d at 20 (rejecting rule that “consent to interception can be inferred from the mere purchase of a service, regardless of circumstances”); *Berry v. Funk*, 146 F.3d 1003, 1011 (D.C. Cir. 1998) (“Without actual notice, consent can only be implied when the surrounding circumstances convincingly show that the party knew about and consented to the interception.”) (quoting *United States v. Lanoue*, 71 F.3d 966, 981 (1st Cir. 1995)) (internal quotations and alterations omitted); *Watkins v. LM Berry & Co.*, 704 F.2d 577, 581 (11th Cir. 1983) (“Consent under title III is not to be cavalierly implied.”).

¹⁹⁵ See Charles Duan, *A Tale of Two Interoperabilities; Or, How Google v. Oracle Could Become Social Media Legislation*, 2021 CARDOZO L. REV. DE-NOVO 246, 252–53 (2021) (noting legislative efforts toward interoperability).

¹⁹⁶ Wiretap Act, 18 U.S.C. § 2511(2)(d).

¹⁹⁷ See *In re Google Inc. Cookie Placement Consumer Priv.*, 806 F.3d 125, 135 (3d Cir. 2015).

¹⁹⁸ See *id.* at 143.

By contrast, in *In re iPhone Application Litigation*, mobile phones were configured to send geolocation information to Apple.¹⁹⁹ Even though that information was sent directly from the phone user to Apple, the district court held the § 2511(2)(d) exception inapplicable, since “[t]he intended communication is between the users’ iPhone and the Wi-fi and cell phone towers,” not Apple’s servers.²⁰⁰ In particular, Apple argued that it was the intended recipient because the phones were designed to transmit geolocation information directly to Apple.²⁰¹ The court rejected this logic on the grounds that it would allow Apple to “manufacture a statutory exception through its own accused conduct.”²⁰²

These cases suggest divergent possible outcomes for the application of the Wiretap Act to message franking. On the one hand, a court could follow *Cookie Placement* and conclude that the franking tag is meant for the messaging platform to review and sign, making the platform the intended recipient. On the other hand, a court following *iPhone* could hold that the intended communication is between the messaging parties and not the platform. The platform’s technical measures to require senders to provide a franking tag might be seen as “manufactur[ing] a statutory exception” through protocol design. Accordingly, it is not certain that this exception would preclude Wiretap Act liability.

Business Use. The other relevant exception is for interception “by a provider of wire or electronic communication service in the ordinary course of its business.”²⁰³ Courts have diverged greatly in their interpretation of this statutory language, particularly with respect to messaging platforms’ automated scanning of messages for purposes of targeted advertising.²⁰⁴ Some narrowly construe

¹⁹⁹ See *In re iPhone Application Litig.*, 844 F. Supp. 2d 1040, 1050–51 (N.D. Cal. 2012).

²⁰⁰ *Id.* at 1062; see also *In re Pharmatrak, Inc. Priv. Litig.*, 329 F.3d 9, 22 (1st Cir. 2003) (holding that a violation of the Wiretap Act can be based on “[s]eparate, but simultaneous and identical, communications” with the interceptor).

²⁰¹ See *iPhone*, 844 F. Supp. 2d at 1062.

²⁰² *Id.*

²⁰³ Wiretap Act, 18 U.S.C. § 2510(5)(a). A separate, related exception relieves an employee of a communications provider from liability for interception “in the normal course of his employment while engaged in any activity which is a necessary incident to the rendition of his service or to the protection of the rights or property of the property of the provider of that service.” *Id.* § 2511(2)(a)(i). It is unclear if this exception applies to service providers themselves. See Boyden, *supra* note 191, at 680.

²⁰⁴ See generally Christopher Batiste-Boykin, *In Re Google Inc.: ECPA*,

“ordinary course of business” to include only interceptions that are “an instrumental part of the transmission” of a message,²⁰⁵ while others more broadly apply the exception to any “customary and routine business practices” of the platform.²⁰⁶

A platform’s interception of a message’s franking tag is almost certainly in the ordinary course of its business under the broader construction, at least if the platform has a customary and routine business practice of content moderation, as most major platforms have. The narrower construction presents a more difficult question, as message franking is not strictly necessary to transmit messages. Nevertheless, message franking as an anti-abuse tool might be analogized to spam filtering or antivirus scanning, technologies that potentially qualify for the exception even under the narrower construction.²⁰⁷

Whether the platform stores any information from the message franking process, such as the franking or verification tags, may affect the analysis under the ordinary-course-of-business exception. Where a platform merely has access to communications content at the time it is being transmitted and not thereafter, courts have held that the platform’s access to the communication is within the ordinary course of its business.²⁰⁸ By contrast, cases dealing with employer recording of telephone calls have held that such recording is not in the ordinary course of business if all calls are recorded.²⁰⁹ As a result, it is possible that the baseline version of message franking, which involves no retention of information, avoids a Wiretap Act violation, while more advanced versions do not.

Consent, and the Ordinary Course of Business in an Automated World, 20 INTELL. PROP. L. BULLETIN 21, 30–34 (2015); Kayla McKinnon, *Nothing Personal, It’s Just Business: How Google’s Course of Business Operates at the Expense of Consumer Privacy*, 33 UIC J. MARSHALL J. INFO. TECH. & PRIV. L. 187, 194–200 (2018); Helen Jazzar, *Bringing an End to the Wiretap Act as Data Privacy Legislation*, 70 CASE W. RES. L. REV. 457, 461–69 (2019).

²⁰⁵ *In re Google Inc. Gmail Litig.*, No. 13-2430, slip op. at 13 (N.D. Cal. Sept. 26, 2013); see *Campbell v. Facebook Inc.*, 77 F. Supp. 3d 836, 844 (N.D. Cal. 2014) (requiring “nexus between . . . the alleged interception and the subscriber’s ultimate business”) (quoting *Gmail*, No. 13-2430, slip op. at 13).

²⁰⁶ See *In re Google, Inc. Priv. Pol’y Litig.*, No. 12-1382, at 19 (N.D. Cal. Dec. 3, 2013).

²⁰⁷ See *Gmail*, No. 13-2430, slip op. at 20 n.4.

²⁰⁸ See *Kirch v. Embarq Mgmt. Co.*, 702 F.3d 1245, 1250 (10th Cir. 2012).

²⁰⁹ See, e.g., *Deal v. Spears*, 980 F.2d 1153, 1158 (8th Cir. 1992).

3. SCA Analysis

There are five points in time when a message and its related franking information are stored, as relevant to the SCA: (1) on the sender's device prior to sending, (2) at the platform while the franking tag is being computed and possibly while the platform is waiting for the recipient to download the message, (3) on the platform after the recipient has downloaded the message, as a backup, (4) on the recipient's device, and (5) on the platform after the recipient has reported abuse.

Section 2701. The general prohibition of the SCA, which covers unauthorized access to a communication service to misuse stored communications, almost certainly does not apply to any of these points in the message franking process, because all of the access to communications is likely authorized and thus not in violation of the statute.²¹⁰ If the platform has not obtained consent from the message sender for the franking process, as described above, then the platform (or the recipient) arguably lacks authorization to use the sender's franking tag.²¹¹ Even so, the statute permits the recipient or the platform itself to authorize access to stored content, making the sender's consent irrelevant.²¹²

The sender of a message might argue, somewhat creatively, that point 1 of sending the message entails an SCA violation to the extent that the sender did not authorize the franking protocol. The argument, akin to the *iPhone* case, would be that the sender's device is a "facility through which an electronic communication service is provided,"²¹³ and that the platform, through its franking-enabled messaging software, accesses messages without authorization before they are sent to construct the franking tag on the sender's device. There are at least three difficulties with this argument. First, it is not clear that an individual user's device can be a "facility" for an "electronic communications service."²¹⁴ Second, at the time the

²¹⁰ See Stored Communications Act, 18 U.S.C. § 2701(a)(1).

²¹¹ See *supra* text accompanying notes 188–95.

²¹² See 18 U.S.C. § 2701(c)(1)–(2).

²¹³ *Id.* § 2701(a)(1).

²¹⁴ *iPhone* considered whether a user device could be considered a "facility" under the SCA. *In re iPhone Application Litig.*, 844 F. Supp. 2d 1040, 1057 (N.D. Cal. 2012). The court held it could not for two reasons. First, such a reading would implausibly mean that "the provider of a communication service could grant access to one's home computer to third parties." *Id.* at 1058 (quoting *Crowley v. CyberSource Corp.*, 166 F. Supp. 2d 1263, 1271 (N.D. Cal. 2001)). Second, treating the user's device as an SCA facility arguably renders the platform a "user" of that facility who can

franking tag is being computed, the message is arguably not in either “temporary, intermediate storage” or “backup protection,” and thus fails to meet the definition of “electronic storage” as the statute requires.²¹⁵ Third and most importantly, even if the sender’s device is a facility of an electronic communications service, the provider of that service (namely, the platform) can authorize access to stored communications on the device.²¹⁶ Accordingly, § 2701 is likely not violated at the time the franking tag is constructed.

Section 2702. The second prohibition of the SCA only concerns the actions of entities providing services “to the public.”²¹⁷ On the assumption that the sender and recipient do not make their devices available to the public, only the platform’s actions at points 2, 3, and 5 above are relevant to this section. Points 2 and 3 only involve disclosure of message information to the intended recipient of the message, which falls cleanly into § 2702’s exceptions.²¹⁸

To the extent that the platform reports the message and related franking information to outside authorities at point 5, the platform presumably has the message recipient’s consent to report the content of the message, which again falls within an exception under § 2702.²¹⁹ However, the platform may be barred from revealing the franking information to law enforcement. Under § 2702(a)(3), a service provider may not “knowingly divulge a record or other information pertaining to a . . . customer of such service . . . to any governmental entity.” Because the franking and verification tags identify the sender of the message, those tags are information

authorize the platform’s access. *See id.* (discussing *Chance v. Ave. A, Inc.*, 165 F. Supp. 2d 1153, 1161 (W.D. Wash. 2001)). A further argument against treating a user device as a facility is that “electronic communications service” is defined as one that provides services “to *users* thereof,” Wiretap Act, 18 U.S.C. § 2510(15) (emphasis added); a single-user device would not seem to fit well within that definition. *See generally* Kerr, *User’s Guide*, *supra* note 66, at 1214–15 & n.47.

²¹⁵ Wiretap Act, 18 U.S.C. § 2510(17); *iPhone*, 844 F. Supp. 2d at 1058–59. This will depend, for example, on whether the message is placed in permanent storage on the sender’s device and whether the franking tag is computed based on that permanently stored copy of the message. *See iPhone*, 844 F. Supp. 2d at 1059 (“Nor do Plaintiffs allege that Defendants accessed the data at a time when the data was only in temporary, intermediate storage.”).

²¹⁶ Stored Communications Act, 18 U.S.C. § 2701(c)(1); *iPhone*, 844 F. Supp. 2d at 1060.

²¹⁷ *See* 18 U.S.C. § 2702(a)(1)–(2).

²¹⁸ *See id.* § 2702(b)(1).

²¹⁹ *See id.* § 2702(b)(3).

pertaining to a customer.²²⁰ Furthermore, the message recipient's consent is irrelevant; consent must originate from "the customer or subscriber" to avoid liability under § 2702(a)(3). Unless the message sender has consented to such disclosure, the platform may only be able to reveal the franking information to law enforcement voluntarily if another statutory exception applies (for example, to protect the service provider, to report emergencies, or to report CSAM).²²¹

4. PRA Analysis

The only candidate for a PR/TT device is the platform's server, at the time it processes a message with a franking tag. A message recipient's own collection of the message is exempt from the statute,²²² and the platform does not act as a PR/TT device upon receiving an abuse report about a communication, because any metadata the platform receives is not collected contemporaneously with the communication itself.²²³

The platform server may qualify as a trap-and-trace device, depending on the construction of "contents" described above.²²⁴ The server captures the sender's identity in order to construct the verification, meaning that the platform server captures information "reasonably likely to identify the source" of the message.²²⁵ However, the server also uses the franking tag, made with a hash of the message contents, to construct the verification tag. If the franking tag is considered contents of the communication, then the platform server falls outside the statutory definition.²²⁶ Notably, a platform's encryption keys may fall within the scope of pen-register interception, according to one court, albeit in a case with an exceptionally unusual procedural posture.²²⁷

²²⁰ As discussed above, some parts of the franking information are arguably content, *see supra* text accompanying notes 181–87, and § 2702(a)(3) does not cover "contents of communications." However, there is metadata in other parts of the franking information. For example, the verification tag includes the sender's identity information.

²²¹ *See* 18 U.S.C. § 2702(c)(3)–(5).

²²² *See* Capitol Recs. Inc. v. Thomas-Rasset, No. 06-1497, slip op. at 8 (D. Minn. June 11, 2009).

²²³ *See* United States v. Fregoso, 60 F.3d 1314, 1321 (8th Cir. 1995).

²²⁴ *See supra* text accompanying notes 181–87.

²²⁵ Pen Register Act, 18 U.S.C. § 3127(4).

²²⁶ *See id.*; *In re* Innovatio IP Ventures, LLC Pat. Litig., 886 F. Supp. 2d 888, 895 (N.D. Ill. 2012).

²²⁷ *See In re* Under Seal, 749 F.3d 276, 292 (4th Cir. 2014). The court primarily held that the service provider had failed to preserve the necessary

Even if the franking tag is not considered contents such that the platform server is a trap-and-trace device, the platform may nevertheless fall into one of several statutory exceptions.²²⁸ First, the sender or recipient may have consented to the server's message franking activities.²²⁹ Message franking might also fall within the exception for "operation, maintenance, and testing of a wire or electronic communication service," akin to the discussion of the business-use exception above.²³⁰

More importantly, the statute exempts a service provider that uses a PR/TT device for "the protection of users of that service from abuse of service or unlawful use of service."²³¹ Message franking, used to support content moderation, would seem to fit cleanly within this exception.

The case law supports this conclusion, though not with complete clarity. While no courts appear to have interpreted the abuse-protection provisions of the PRA, several state courts have construed analogous provisions of state-law equivalents to the federal statute in the context of telephone caller identification services. The South Carolina Supreme Court applied that exception, holding that the caller ID service "is designed to protect the utility's subscribers from abusive or unlawful telephone calls."²³² By contrast, the Pennsylvania appellate and supreme courts did not address the abuse-protection exceptions in holding caller ID services to be illegal trap-and-trace devices.²³³ Although it is not clear why, the appellate decision expressed general skepticism about the service's likelihood of preventing abuse, finding it "conceivable that Caller*ID is just as likely to encourage criminal or annoying behavior as it would to discourage such conduct."²³⁴ These decisions suggest that some courts might simply accept that content moderation is an abuse-protection objective that exempts platforms from PR/TT device regulation, while other courts might take a harder look at the platform's specific content moderation policies and practices to

arguments. *See id.* at 293.

²²⁸ *See* 18 U.S.C. § 3121(b).

²²⁹ *See supra* text accompanying notes 188–95.

²³⁰ *See* 18 U.S.C. § 3121(b)(1); *supra* text accompanying notes 203–09.

²³¹ 18 U.S.C. § 3121(b)(1); *see also id.* § 3121(b)(2) (providing exception for use of PR/TT devices "to record the fact that a wire or electronic communication was initiated or completed in order to protect . . . a user . . . from fraudulent, unlawful or abusive use of service"); *Barasch v. Pa. Pub. Util. Comm'n*, 409 S.E.2d 775, 777 (S.C. 1991).

²³³ *See Barasch v. Pa. Pub. Util. Comm'n*, 576 A.2d 79, 85–86 (Pa. Commw. Ct. 1990), *aff'd sub nom. Barasch v. Bell Tel. Co. of Pa.*, 605 A.2d 1198 (Pa. 1992).

²³⁴ *Id.* at 90.

decide whether the exception applies.

5. CFAA Analysis

The question to be answered under the CFAA is whether the platform, in the course of the message franking protocol, accesses a protected computer in violation of the statute.²³⁵ Since the platform has authorization to access its own servers, the sender's device and the receiver's device are the two primary computers to be considered. Per Jonathon W. Penney and Bruce Schneier, the overall encrypted messaging network could further be a protected computer.²³⁶

With respect to the sender's device, the argument would be that the platform, by adding message franking features to the end-to-end encrypted messaging software the sender uses, accesses the sender's unencrypted message without authorization in order to generate the franking tag.²³⁷ In a sense, the sender would argue that the message franking features are a form of spyware, working around the sender's expectation of privacy through end-to-end encryption.²³⁸

If the sender explicitly authorizes the franking tag generation (say, by accepting the platform's terms of service), then there is likely no violation of the CFAA. And even absent explicit authorization, a court might find implicit authorization based on the operation of the franking protocol.²³⁹

For the recipient's device, the platform installs software that, upon receipt of a message, verifies the franking tag against the sender's random key and then blocks or otherwise affects display of the message if the verification fails. If the recipient consents to the

²³⁵ See Computer Fraud and Abuse Act, 18 U.S.C. § 1030(a)(2)(C).

²³⁶ See Penney & Schneier, *supra* note 116, at 488–90.

²³⁷ See Kerr & Schneier, *supra* note 27, at 1007–10 (describing techniques for accessing on-device plaintexts to circumvent encryption).

²³⁸ *Cf. id.* at 1009 (describing use of “government malware” to obtain IP addresses of encryption users); James Grimmelmann, *Spyware vs. Spyware: Software Conflicts and User Autonomy*, 16 OHIO ST. TECH. L.J. 25, 58–59 (2020) [hereinafter Grimmelmann, *Spyware*] (questioning determinations of user consent when two pieces of software operate to contrary ends).

²³⁹ See *supra* text accompanying notes 188–95 (describing how sender's inclusion of a franking tag might constitute implicit consent); *hiQ Labs, Inc. v. LinkedIn Corp.*, 31 F.4th 1180, 1197–98 (9th Cir. 2022) (finding no lack of consent based on website's actions of making information publicly accessible); Grimmelmann, *Consenting*, *supra* note 123, at 1508 (considering situations where a party “has done something . . . that manifests her factual consent”); *cf. EF Cultural Travel BV v. Zefer Corp.*, 318 F.3d 58, 63 (1st Cir. 2003) (“[L]ack of authorization may be implicit, rather than explicit.”).

verification process, then the platform's access to the recipient's device is authorized and no violation of the CFAA occurs. Even if the recipient fails to authorize the verification process, no violation occurs because there is no actionable harm. The platform does not obtain any information or thing of value from the recipient's device.²⁴⁰ The platform's software arguably "causes damage" to the recipient's device by blocking the "availability of data" on that device, which might violate 18 U.S.C. § 1030(a)(5).²⁴¹ But a violation of that part of the CFAA requires access "without authorization," and by virtue of installing the platform's software voluntarily, the recipient provided the necessary authorization.²⁴²

Yet, even though the platform does not violate the CFAA for each device individually, there is a not-completely-implausible argument that the platform violates the CFAA by trespassing upon the network as a whole. Under Penney and Schneier's network-trespass theory, the "computer," for purposes of the CFAA, is the entire network of messaging participants, including the platform's servers and the devices of message senders and recipients.²⁴³ The platform obviously is authorized to access its own network, but it is not authorized to access every piece of data thereon. After all, the whole point of end-to-end encryption is that the platform has no access to the communicated information.²⁴⁴ Whether the platform exceeds authorized access then depends on what information the platform's users intend to shield from the platform with their use of encryption.

There is a good argument that one piece of information that the participants intend to shield is the fact that a specific sender sent a specific piece of content. In an ordinary encrypted messaging system, the platform cannot prove that a certain user sent a particular encrypted message.²⁴⁵ Even if the recipient of a message reveals the message to the platform, the platform cannot be sure who sent the message unless the sender chose to include a digital signature or

²⁴⁰ See 18 U.S.C. § 1030(a)(2)–(3).

²⁴¹ See *id.* § 1030(a)(5)(A), (e)(8).

²⁴² Unlike other violations of the CFAA, damage to a computer is only actionable based on access "without authorization." *Id.* § 1030(a)(5); *Int'l Airport Ctrs., LLC v. Citrin*, 440 F.3d 418, 420 (7th Cir. 2006) (noting that 18 U.S.C. § 1030(a)(5) does not prohibit damage from "exceeding authorized access").

²⁴³ See Penney & Schneier, *supra* note 116, at 490–93.

²⁴⁴ See *id.* at 494.

²⁴⁵ See Grubbs et al., *supra* note 162, at 67 ("But end-to-end confidentiality means that Facebook must rely on users sending examples of malicious messages. How can the provider know that the reported message was the one sent?").

some other authenticating information with the message.

With message franking, the platform receives new information, previously inaccessible due to encryption, about who sent a message. The spyware argument that failed with respect to the sender's device alone potentially succeeds under the network-trespass theory, because the sender can argue that the platform's software programs on the sender's and recipient's devices *together* are the "spyware" that let the platform exceed authorized access to information about the message sender's identity. Accordingly, the platform exceeds its authorized access to the encrypted messaging channel to obtain otherwise-inaccessible information, meeting all the elements of § 1030(a)(2)(C).

6. CALEA Analysis

With respect to CALEA, the question is whether a messaging platform implementing message franking would be required to implement that technology in some manner to enable interception by law enforcement.²⁴⁶ There are three candidate communications to which the statute might apply: (1) the transmission of an encrypted message across the platform, (2) the transmission of associated franking information with a message, and (3) a recipient's report of an abusive message sent to the platform. The first plainly falls within CALEA's encryption exception because, by definition, an end-to-end encrypted platform denies the platform access to the encryption keys.²⁴⁷ And assuming that law enforcement has adequate authorization to demand such interception, no technical capability is required for the third, since the platform can simply transmit to law enforcement anything it receives from the recipient's report.

However, CALEA may require covered platforms to build in interception capabilities for the message franking and verification tags as they are sent across the platform. The threshold question is whether the platform is a "telecommunications carrier."²⁴⁸ Most online platforms will fall under CALEA's expansive definitions of "electronic messaging service" or "information service," which do not qualify as telecommunications carriers.²⁴⁹ However, a synchronous voice-based messaging platform might be deemed sufficiently a "replacement for a substantial portion of the local telephone service," such that the platform could be deemed a

²⁴⁶ See Communications Assistance for Law Enforcement Act (CALEA) § 103(a)(1), 47 U.S.C. § 1002(a)(1).

²⁴⁷ See *id.* § 103(b)(3).

²⁴⁸ See *id.* § 103(a).

²⁴⁹ See *id.* § 102(4), (6), (8)(C)(i).

“telecommunications carrier” for purposes of the statute.²⁵⁰

If the platform qualifies as a telecommunications carrier, then CALEA’s exceptions likely do not apply.²⁵¹ The technological-implementation limitation of section 103(b)(1) probably does not apply so long as law enforcement does not seek to require message franking or demand that specific franking software be used.²⁵² The information-service exception does not apply if the platform has been determined to be a telecommunications carrier.²⁵³ The encryption exception potentially does not apply because the “encryption was provided by the carrier” (the platform’s software that implements franking) and “the carrier possesses the information necessary to decrypt the communication” (the platform’s secret key used to encrypt the verification tag).²⁵⁴

To the extent CALEA requires interception capabilities of platforms with message franking, what capabilities must be included? Likely the franking and verification tags would have to be retained and delivered to the government upon appropriate authorization, as that information could be deemed “call-identifying information that is reasonably available to the carrier.”²⁵⁵ However, by design, those tags provide virtually no information without access to the unencrypted message content. Law enforcement could also obtain the sender information context that the platform uses to construct the verification tag; that information could overcome anonymity guarantees on platforms allowing for anonymous sending of messages.

The most concerning possibility would be that, should CALEA apply to a platform with message franking, then the government could initiate a standard-setting process at the FCC and propose designs for message franking that would weaken the end-to-end encryption guarantees of the messages themselves. It is unlikely that such a proposal would succeed, given that it would be an end run around the statute’s encryption exception and would also violate CALEA’s warning that the statute “does not authorize any law enforcement agency or officer . . . to require any specific design of

²⁵⁰ *Id.* § 102(8)(B)(ii); *see* *Am. Council on Educ. v. FCC*, 451 F.3d 226, 232–33 (D.C. Cir. 2006).

²⁵¹ *See* CALEA § 103(b).

²⁵² *See id.* § 103(b)(1).

²⁵³ *See id.* § 103(b)(2). To be sure, information services that the platform provides would still be exempt from CALEA’s requirements. *See Am. Council on Educ.*, 451 F.3d at 233 (holding that “‘telecommunications carrier’ and ‘information services’ are not mutually exclusive terms”).

²⁵⁴ CALEA § 103(b)(3).

²⁵⁵ *Id.* § 103(a)(2).

equipment, facilities, services, features, or system configurations.”²⁵⁶ Nevertheless, the possibility of a protracted technical standards battle points to a risk that implementers of message franking may face, should they implement that technology on voice-like services potentially within the ambit of CALEA.

7. POCA Analysis

Since the baseline message franking protocol does not involve hashes or other elements of content,²⁵⁷ the relevant provisions of POCA are the mandatory reporting requirement of § 2258A and the immunity of § 2258B. Regarding the former, the statute comes into play only upon a platform “obtaining actual knowledge” of transmission of CSAM.²⁵⁸ Since the process of generating and verifying franking tags retains the privacy of the encrypted communication,²⁵⁹ any actual knowledge would only arise after a message recipient makes a report to the platform. The difference that message franking makes is that the platform is able, though not required, to include the franking and verification tags as cryptographically verifiable “[i]nformation relating to the identity of any individual who appears to have violated” the relevant CSAM laws.²⁶⁰

Regarding immunity, the key point is that § 2258B covers only acts of disclosure, not acts of searching or reading content.²⁶¹ As a result, this immunity does not overcome potential liability under most of the communication privacy laws discussed above, as those laws focus on acts of interception or unauthorized access prior to any act of reporting. The SCA does prohibit divulging protected communications,²⁶² so the § 2258B immunity would apply to that

²⁵⁶ *Id.* § 103(b)(1)(A).

²⁵⁷ Experimental message franking protocols may use a hash database—for example allowing the platform to uncover a message sender’s identity only when the message matches content in a hash database. *See* James Bartusek, Sanjam Garg, Abhishek Jain & Guru-Vamsi Policharla, *End-to-End Secure Messaging with Traceability Only for Illegal Content*, 42 ANN. INT’L CONF. ON THEORY & APPLICATIONS CRYPTOGRAPHIC TECHS. 35, 37–42 (2023) (described *infra* note 281). For such systems, 18 U.S.C. § 2258C may be relevant.

²⁵⁸ PROTECT Our Children Act of 2008, 18 U.S.C. § 2258A(a)(1)(A).

²⁵⁹ *See, e.g.*, Grubbs et al., *supra* note 162, at 75.

²⁶⁰ 18 U.S.C. § 2258A(b)(1).

²⁶¹ *See id.* § 2258B(a) (limiting liability for acts “arising from the performance of the reporting or preservation responsibilities” of § 2258A); *United States v. Stevenson*, 727 F.3d 826, 830 (8th Cir. 2013).

²⁶² *See* Stored Communications Act, 18 U.S.C. § 2702(a)(1).

statute, but the SCA also contains internal exceptions for reporting CSAM.²⁶³ As a result, § 2258B likely has little effect on platforms' liability for content moderation activities under the communication privacy laws.

B. Forward Tracing

Our next content moderation technology is “forward tracing.” Message franking can help platforms identify senders of reported messages, so that the platform can respond to improper content. However, content moderators are often concerned not just with single messages that are sent once, but with messages that are forwarded, potentially many times. When a message recipient reports the message to the platform for moderation, the platform may wish to know not just the immediate sender of that message but also all of the prior senders in the forwarding chain, so that it can identify the origin of the message.

Forward-tracing protocols overcome this limitation of message franking, helping platforms determine the originator of a given message. To do this, the protocol must keep track of information about the message's forwarding path, in addition to its contents and regular metadata. In an E2EE system, the challenge for such protocols is how to track this forwarding information without defeating the privacy expectations that E2EE provides.

²⁶³ See *id.* § 2702(b)(6), (c)(5). The Wiretap Act, in addition to its prohibition on interception, contains a separate prohibition on divulging contents of communications to third parties. See Wiretap Act, 18 U.S.C. § 2511(3)(a). A violation of this prohibition would be immunized under § 2258B when a platform makes a CSAM report, but § 2258B would not immunize any underlying act of interception under (1)(a). In any event, the Wiretap Act divulgence prohibition itself also contains an exception for reporting apparent crimes to law enforcement. See *id.* § 2511(3)(b)(iv).

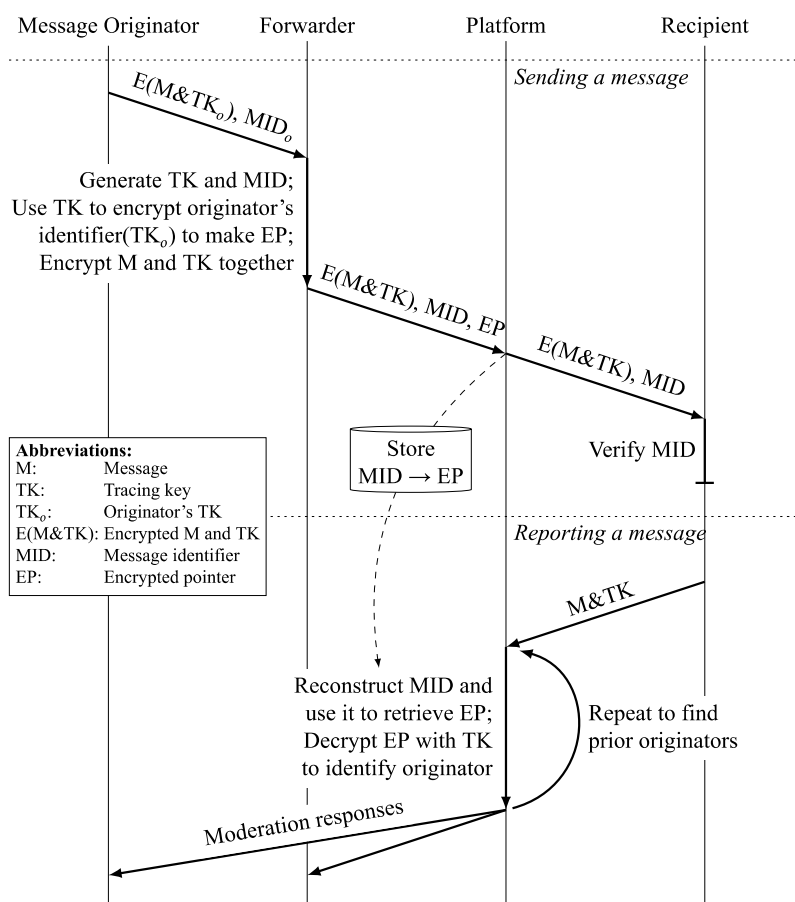


Figure 2: Diagram of communications performed in traceback.

1. Technical Overview

Researchers Nirvan Tyagi, Ian Miers, and Thomas Ristenpart published the basic proposal for message traceback in 2019.²⁶⁴ The protocol begins much like message franking: to send a message, the sender first produces a “message identifier” based on a cryptographic hash of the message content and a randomly generated secret key called the “tracing key.”²⁶⁵

Next, the sender constructs an “encrypted pointer” that identifies the forwarding origin of the message. If the sender is forwarding

²⁶⁴ See Nirvan Tyagi, Ian Miers & Thomas Ristenpart, *Traceback for End-to-End Encrypted Messaging*, PROC. ACM SIGSAC CONF. ON COMPUT. & COMM’NS SEC. 413 (2019).

²⁶⁵ See *id.* at 416.

another message, then the encrypted pointer is the prior message's tracing key, encrypted with the newly generated tracing key.²⁶⁶ As a result, the encrypted pointer chains the forwarding path together: the tracing key of the last message in the forwarding chain can unlock the previous message's tracing key, that tracing key can unlock its predecessor, and so on.²⁶⁷ If the sender has composed an original, non-forwarded message, then the sender constructs and encrypts a nonexistent tracing key, effectively cutting off the forwarding chain.²⁶⁸

The sender then encrypts the message and the tracing key together using the message recipient's public key, sending the resulting ciphertext along with the message identifier and the encrypted pointer. Upon receipt, the platform stores the message identifier and encrypted pointer, and then delivers the message ciphertext and message identifier to the recipient. The recipient decrypts the message and verifies that the message identifier was constructed correctly before displaying the message.²⁶⁹

At this point, neither the platform nor the message recipient can determine the message forwarding chain. The message identifier is a cryptographic hash based in part on a random key, so it cannot be used to identify related messages. The encrypted pointer contains the prior message's tracing key, meaning that it can help to identify the prior message in the chain if decrypted. But the platform lacks the tracing key to decrypt the pointer, and the recipient does not receive the pointer.²⁷⁰ As a result, neither the platform nor the recipient can determine whether the message is the sender's original content or forwarded from someone else.²⁷¹

When a message recipient wishes to report an abusive message, the recipient sends the platform the unencrypted message content and tracing key for the message. With these pieces of information, the platform can reconstruct the message identifier and thus find the associated encrypted pointer. The platform can then decrypt the encrypted pointer with the tracing key to discover the prior message's tracing key, reconstruct the prior message's identifier, find that prior message's encrypted pointer, and so on until the entire

²⁶⁶ *See id.* at 417.

²⁶⁷ *See id.*

²⁶⁸ *See id.*

²⁶⁹ *See id.* Note that neither the platform nor the recipient verifies the encrypted pointer. Doing so would have limited value, because a determined message forwarder can always make it appear that a message has not been forwarded, for example by copying and pasting content into a new, apparently unforwarded message. *See id.* at 415.

²⁷⁰ *See* Tyagi et al., *supra* note 264, at 418.

²⁷¹ *See id.* at 417–18.

forwarding chain has been revealed.

An alternate approach, proposed by researchers Charlotte Peale, Saba Eskandarian, and Dan Boneh, is to track just the original sender of a message rather than the entire forwarding chain.²⁷² Although it gives the platform less information to work with in responding to problematic content, this “source-tracking” approach requires no platform-side storage, and it provides message senders a greater degree of privacy than the traceback protocol.²⁷³

Source tracking can be understood as a variant of message franking. For all messages, the platform produces a signature analogous to the verification tag of message franking and based on the sender’s identity and a message “commitment” (analogous to the franking tag).²⁷⁴ The recipient of an original, unforwarded message receives and verifies the platform-produced signature.²⁷⁵ When forwarding a message, however, the sender includes the platform signature inside the message to be encrypted.²⁷⁶ The platform, not knowing whether the message ciphertext contains a signature inside, generates a new signature, but the recipient of the forwarded message discards this generated signature and verifies only the signature found inside the message.²⁷⁷ As a result, all forwards of a message will internally contain the first platform signature generated for the message, which the platform can use to identify the original sender when the message is reported.²⁷⁸

Further research improves on the basic traceback and source-tracing schemes. Tyagi and colleagues also propose a mechanism for tracing further downstream recipients who received a reported message (perhaps useful, for example, where the abusive message is a phishing scam so the platform can advise those other recipients) by having the platform store further information for downstream tracing.²⁷⁹ Others extend the traceback protocol to sender-anonymous platforms.²⁸⁰ James Bartusek, Sanjam Garg, Abhishek

²⁷² See Charlotte Peale, Saba Eskandarian & Dan Boneh, *Secure Complaint-Enabled Source-Tracking for Encrypted Messaging*, PROC. ACM SIGSAC CONF. ON COMPUT. & COMM’NS SEC. 1484 (2021); Issa et al., *supra* note 176 (source tracking in sealed-sender message systems).

²⁷³ See Peale et al., *supra* note 272, at 1485.

²⁷⁴ See *id.* at 1491.

²⁷⁵ See *id.* at 1491 fig.3.

²⁷⁶ See *id.*

²⁷⁷ See *id.* at 1491. In particular, Peale and colleagues specify giving the server a nonsense message commitment for a forwarded message, ensuring that the platform’s generated signature is useless. See *id.*

²⁷⁸ See Peale et al., *supra* note 272, at 1491.

²⁷⁹ See Tyagi et al., *supra* note 264, at 420–21.

²⁸⁰ See Erin Kenney, Qiang Tang & Chase Wu, *Anonymous Traceback for*

Jain, and Guru-Vamsi Policharla propose further mechanisms that limit the platform's ability to decrypt traceback records, and thereby discover the originator of a message, unless the message matches a database of illicit content.²⁸¹ Linsheng Liu, Daniel S. Roche, Austin Theriault, and Arkady Yerukhimovich similarly propose a system that limits platforms' ability to uncover the original sender of a message, requiring first that a threshold number of users report the message before the sender's identity can be decrypted.²⁸²

2. Wiretap Act Analysis

For purposes of the Wiretap Act, forward-tracing protocols are largely identical to message-franking protocols, so the prior analysis applies with equal force.²⁸³ The message identifier for traceback protocols and the message commitment for source tracking are both hashes computed on the message plaintext, so they are content to the same extent that a franking tag is content.²⁸⁴ The consent exception focuses on the parties to the communication at the time of transmission, so the consent of the original message sender is probably irrelevant and the exception turns on the consent of the forwarder and forward recipient.²⁸⁵ The intended-recipient exception would also be analyzed in the same way as it was for message franking.²⁸⁶

The business-use exception does not necessarily apply in the same way it does for message franking, because the platform's

End-to-End Encryption, 27 EUR. SYMP. ON RSCH. COMPUT. SEC. 42 (2022).

²⁸¹ See Bartusek et al., *supra* note 257, at 37–38. At a high level, the sender of a message uses the message content and a specially generated “set pre-constrained encryption” key to encrypt identity traceback information. See *id.* at 42. The corresponding private key, held by the platform, is designed based on the database of illicit content such that it can only decrypt the traceback information if the message content was contained in that database. See *id.* at 40.

²⁸² See Linsheng Liu, Daniel S. Roche, Austin Theriault & Arkady Yerukhimovich, *Fighting Fake News in Encrypted Messaging with the Fuzzy Anonymous Complaint Tally System (FACTS)*, NETWORK & DISTRIBUTED SYS. SEC. SYMP. (2022), <https://www.ndss-symposium.org/wp-content/uploads/2022-109-paper.pdf> [<https://perma.cc/WSS9-BY2Q>].

²⁸³ See *supra* Part III.A.2.

²⁸⁴ See *supra* text accompanying notes 181–87.

²⁸⁵ See *supra* text accompanying notes 188–95; Wiretap Act, 18 U.S.C. § 2511(2)(d) (considering whether “one of the parties to the communication has given prior consent”).

²⁸⁶ See *supra* text accompanying notes 196–202.

purpose for forward tracing differs from that for message franking.²⁸⁷ Unlike abuse prevention tools like spam filtering or antivirus scanning that courts have suggested might satisfy the exception,²⁸⁸ forward tracing provides information external to an individual message transaction, namely information about who else sent or received a message with identical content. A court adopting a narrower reading of the business use exception, then, may see interception of forward tracing data as lacking a sufficient nexus with the platform's business purposes of message transmission.²⁸⁹ This is especially so for traceback, since the platform permanently stores every message's identifier regardless of whether the platform suspects wrongdoing justifying interception.²⁹⁰ To the extent that message franking presented several difficult analytical questions under the Wiretap Act,²⁹¹ forward tracing enhances the difficulty of those questions.

3. SCA Analysis

As with message franking, there is likely no violation of the general unauthorized-access provision of the SCA, because all data access during a forward-tracing protocol can be authorized by the platform implementing the protocol.²⁹² Therefore, we focus here on § 2702's rules relating to service providers' disclosure of stored communications.

For forward tracing, information is divulged to third parties at two possible points: (1) where traceback is used and a user has reported an illicit message, the platform may provide message identifiers or encrypted pointers to law enforcement or others involved in content moderation, and (2) where sender tracking is used, the platform's signature data is passed along the forwarding chain.

In the first case, where the platform discloses information to law enforcement or others, the disclosure can meet all the requirements of § 2702(a)(1), though it is unlikely. The platform is an electronic communication service available to the public. The message and the list of users are knowingly divulged. The tracing information proving

²⁸⁷ Cf. *supra* text accompanying notes 203–09.

²⁸⁸ See *In re Google Inc. Gmail Litig.*, No. 13-2430, slip op. at 20 n.4 (N.D. Cal. Sept. 26, 2013).

²⁸⁹ *Campbell v. Facebook Inc.*, 77 F. Supp. 3d 836, 844 (N.D. Cal. 2014); *Gmail*, No. 13-2430, slip op. at 13.

²⁹⁰ See *supra* text accompanying notes 208–09.

²⁹¹ See *supra* Part III.A.2.

²⁹² See Stored Communications Act, 18 U.S.C. § 2701(a)(1); *supra* text accompanying notes 210–16.

the list of users is arguably “contents of a communication,” because the message identifier of the traceback protocol is cryptographically derived from the message content.²⁹³ The more difficult question is whether the traceback data is in electronic storage. It certainly is not in temporary, intermediate storage since the platform stores it after the communication transaction has concluded.²⁹⁴ And it is difficult to see how the traceback data serves as “backup protection,” even under *Theofel*, because the platform does not store that data to help users recover lost message identifiers or encrypted pointers.²⁹⁵ Nevertheless, a court might be persuaded that the platform’s retention of traceback data for later content moderation accountability is a kind of “backup” that meets the statutory definition.²⁹⁶

The remote computing service provision of § 2702(a)(2) probably does not apply to forward-tracing platforms because the traceback or source-tracking data is not provided to the platform solely for storage and computer processing purposes (they are intended to be forwarded to the message recipient for verification purposes).²⁹⁷ However, the prohibition on disclosing customer information under § 2702(a)(3) might apply if the platform reveals traceback data to law enforcement, as that data identifies message senders so it is “information pertaining to a subscriber.”²⁹⁸

Even if § 2702(a)(1) or (a)(3) applies as above, the platform may satisfy one or more of its exceptions. Message senders and recipients may have consented to the disclosure of traceback information as part of the platform’s terms of service.²⁹⁹ Unlike message franking, however,³⁰⁰ the consent of the person reporting an improper message does not exempt disclosure of the entire forwarding chain, because

²⁹³ See *supra* text accompanying notes 181–87.

²⁹⁴ Compare Wiretap Act, 18 U.S.C. § 2510(17)(A), with Tyagi et al., *supra* note 264, at 416 fig.2.

²⁹⁵ See *Theofel v. Farey-Jones*, 359 F.3d 1066, 1076 (9th Cir. 2004) (“But the mere fact that a copy *could* serve as a backup does not mean it is stored for that purpose.”); *Republic of the Gam. v. Facebook, Inc.*, 567 F. Supp. 3d 291, 305–06 (D.D.C. 2021) (“Facebook claims it kept the instant records as part of an autopsy of its role in the Rohingya genocide. . . . While admirable, that is storage for self-reflection, not for backup.”); *Quon v. Arch Wireless Operating Co.*, 529 F.3d 892, 902 (9th Cir. 2008) (holding that text messaging service Arch Wireless held already-delivered text messages for “backup protection” even though it was “not clear for whom Arch Wireless ‘archived’ the text messages”).

²⁹⁷ See 18 U.S.C. § 2702(a)(2)(B).

²⁹⁸ *Id.* § 2702(a)(3).

²⁹⁹ See *id.* § 2702(b)(3), (c)(2).

³⁰⁰ See *supra* text accompanying note 219.

other messages in the chain may not have involved the recipient at all.³⁰¹ The platform could also argue that forward tracing, as a mechanism for accountability and abuse prevention, is a necessary incident of running an encrypted messaging service and protects the platform's rights and property.³⁰² Finally, if the disclosure is made to law enforcement to prevent serious harm or report child exploitation, then further exceptions may apply.³⁰³

Regarding the platform's passing of sender tracking data to forward recipients, an argument might be that the sender's message commitment, as part of the platform's sender-tracking signature, has been divulged not just to the original message recipient but also to all third parties to whom the message is later forwarded. This argument fails both because the message recipient consents to the inclusion of the commitment and signature in future forwards, and it is the message forwarders, not the platform itself, who are divulging those data elements.

4. *CALEA Analysis*

Under CALEA, a law enforcement agency could require a platform implementing a forward tracing protocol to capture message identifiers, encrypted pointers, or other information generated and transmitted to the platform in the course of the protocol.³⁰⁴ As with message franking, this information is theoretically meaningless alone, but in combination with message plaintexts that law enforcement might obtain through legal or investigative means, the intercepted forward-tracing information could give law enforcement access to identities of participants in a forwarded conversation chain.

For CALEA to apply, the threshold question is whether the messaging service implementing forward tracing is a telecommunications carrier under the statute, as opposed to an information service. Yet the service of forwarding messages, as required for any forward-tracing protocol, would not seem to be a "replacement" for traditional telephone services, as the statute requires of a telecommunications carrier.³⁰⁵ As a result, a platform

³⁰¹ For example, say that A sends a message to B who forwards it to C, and C reports the message to the platform. C's consent cannot excuse disclosing traceback information between A and B.

³⁰² See 18 U.S.C. § 2702(b)(5), (c)(3).

³⁰³ See *id.* § 2702(b)(6)–(8), (c)(4)–(5).

³⁰⁴ See Communications Assistance for Law Enforcement Act (CALEA) § 103(a), 47 U.S.C. § 1002.

³⁰⁵ *Id.* § 102(8)(B)(ii); see *Am. Council on Educ. v. FCC*, 451 F.3d 226, 232 (D.C. Cir. 2006).

would have a strong argument that the forward-tracing protocol is part of an information service and thus not susceptible to the capabilities requirements of CALEA.

If the forward-tracing platform is considered a telecommunications carrier, then the next question is whether CALEA's encryption exception exempts the platform. It probably does not. For a source-tracking protocol, the platform has the key used to encrypt the signature identifying the message source.³⁰⁶ For a traceback-based protocol, however, the platform does not possess the tracing key used to produce the encrypted pointer.³⁰⁷ But on the assumption that law enforcement has obtained the tracing key from a message recipient, all that is necessary is for the platform to deliver the encrypted pointer for law enforcement to decrypt.

To the extent that the platform is deemed a telecommunications carrier, CALEA effectively places law enforcement in a privileged position above the platform itself. Like the platform, law enforcement has the ability to uncover the original sender, chain of forwarders, or the entire tree of message recipients (depending on the forward tracing protocol the platform implements), so long as one plaintext message is revealed. Unlike the platform, however, law enforcement enjoys a range of compulsory legal and investigative powers to cause disclosure of that one plaintext message that unlocks the intercepted forward tracing information.

5. PRA Analysis

The analysis of the PRA for forward tracing schemes is essentially identical to the analysis for message franking. The platform server is a trap-and-trace device as long as the message hashes (the message identifier for traceback or the message commitment for source tracking) are not content, and the statutory exceptions for user consent, service operations, and abuse protection will likely apply.³⁰⁸

6. CFAA Analysis

The forward-tracing analysis under the CFAA similarly mirrors that for message franking. With respect to any individual user's device, the platform has authorization to generate tracing keys or

³⁰⁶ See Peale et al., *supra* note 272, at 1491; CALEA § 103(b)(3) (stating that encryption exception does not apply if "the carrier possesses the information necessary to decrypt the communication").

³⁰⁷ See Tyagi et al., *supra* note 264, at 417.

³⁰⁸ See *supra* Part III.A.4.

message identifiers, either based on explicit user consent or by implicit consent in order for a message to be verifiable upon transmission.³⁰⁹ If the encrypted messaging network is considered a single “computer” for purposes of the CFAA, then the original sender of a message could argue, analogously to message franking, that the platform’s software for forward tracing circumvents the privacy guarantees of end-to-end encryption and therefore exceeds authorized access.³¹⁰

7. POCA Analysis

A reporting requirement under § 2258A arises only after a user report since all information prior to then is encrypted and therefore cannot give rise to actual knowledge.³¹¹ And the liability shield of § 2258B only affects potential liability under the SCA, which itself already provides an exception for divulging CSAM-related materials to NCMEC.³¹²

C. Server-Side Automated Content Scanning

Automated content scanning, where predefined algorithms of varying complexity sort out the good content from the bad, is a widely debated technique for content moderation.³¹³ But putting aside the debate over automated moderation’s effectiveness, end-to-end encryption raises a more basic question: if encryption prevents a platform from reading content, then can the platform apply algorithmic filtering in the first place? Surprisingly, it can.

1. Technical Background

This automated server-side filtering depends on a class of algorithms known as homomorphic encryption.³¹⁴ These systems have the property that computations done on the encrypted ciphertext

³⁰⁹ See *supra* Part III.A.5.

³¹⁰ See *supra* text accompanying notes 243–45.

³¹¹ See PROTECT Our Children Act of 2008, 18 U.S.C. § 2258A(a)(1)(A).

³¹² See *supra* Part III.A.7.

³¹³ See, e.g., Tarleton Gillespie, *Content Moderation, AI, and the Question of Scale*, BIG DATA & SOCIETY (2020); Daphne Keller, *Internet Platforms: Observations on Speech, Danger, and Money*, HOOVER INST. 5-8 (2018), <https://www.hoover.org/research/internet-platforms-observations-speech-danger-and-money> [https://perma.cc/4CKV-LD7C].

³¹⁴ See, e.g., NAT’L ACAD. OF SCI., ENG’G & MED., *supra* note 24, at 31; WONG, *supra* note 12, § 15.2.

will predictably operate on the underlying unencrypted plaintext, so that the results of the computation can be retrieved once the message is decrypted.³¹⁵ The ROT-13 encryption cipher described previously exemplifies this homomorphic property for a number of computations such as text reversal. Consider the following series of operations:³¹⁶

$$\begin{aligned} \text{"GOHANGASALAMI"} &\xrightarrow{\text{Encrypt}_{\text{ROT-13}}} \text{"TBUNATNFNYNZV"} \\ \text{"TBUNATNFNYNZV"} &\xrightarrow{\text{Compute}} \text{"VZNYNFNTANUBT"} \\ \text{"VZNYNFNTANUBT"} &\xrightarrow{\text{Decrypt}_{\text{ROT-13}}} \text{"IMALASAGNAHOG"} \end{aligned}$$

Importantly, the fact that a platform can perform computations on homomorphically encrypted content does not mean that the platform gains any insight into the nature of the content. The content after computation appears just as scrambled and meaningless as the original encrypted content, and the result of the computation can only be perceived after decryption.³¹⁷ For this reason, homomorphic encryption enables a content-scanning system to act upon messages without requiring users to expose the plaintext of those messages.

However, the actions available to the platform-based content scanner are quite limited. The scanner cannot itself determine whether content is flagged or problematic, since the results of the scanner's computations are buried within the content's encryption. Instead, the scanner can only modify the content that the recipient will see, because modifications to content are simply the results of computations on that content. For example, the scanner can attach a flag to content, or theoretically even blur or black out undesirable images. Message recipients would learn of the platform's modifications upon decrypting the messages,³¹⁸ but the platform would not thereby learn whose messages were flagged or modified.³¹⁹

³¹⁵ See WONG, *supra* note 12, § 15.2.

³¹⁶ With apologies to John Agee, who devised this palindrome. JON AGE, GO HANG A SALAMI! I'M A LASAGNA HOG! AND OTHER PALINDROMES (1994).

³¹⁷ See WONG, *supra* note 12, § 15.2 ("The important idea here is that the service never learns about your values and always deals with ciphertexts.").

³¹⁸ If the platform could modify content without the recipient detecting the modification, then the platform could effectively perform "man-in-the-middle" attacks, defeating the authenticity guarantee that end-to-end encrypted systems typically provide.

³¹⁹ If the client device is configured to report the outcome of the server-side scan back to the platform or to a third party, then such a system would

Homomorphic encryption is not practically usable for content moderation, because current algorithms are too slow to be used at the scale of large messaging platforms.³²⁰ Nevertheless, several researchers have proposed server-side scanning systems for content moderation using homomorphic encryption.³²¹ Furthermore, related technologies, such as secure multi-party communication and functional encryption, similarly may allow a platform to perform computations on a message without giving the platform access to the message's content.³²² Future advances in cryptographic algorithms may thus create greater opportunities for platforms to moderate messages without reading them.

2. Wiretap Act Analysis

For platform-based automated content scanning, the relevant point of interception occurs when the platform's server receives the homomorphically encrypted message and performs computations on it. The message is an electronic communication under the statute, the server is a device, and the platform acts intentionally by using automated scanning, so there is a *prima facie* violation of the Wiretap Act if the encrypted message qualifies as "contents" and the platform "intercepts" it.³²³

Bruce E. Boyden has argued that no interception should be found based on purely automated message processing—for example, to append advertisements to the message.³²⁴ However, in two Wiretap Act cases involving automated message processing, both postdating Boyden's article, the platform defendants did not raise this argument

be classified as client-side scanning because the client device is performing automatic processing on the communication, albeit with assistance from the server-side scanner. *See infra* Part III.D.1.

³²⁰ *See* WONG, *supra* note 12, § 15.2.5 ("At the time of this writing (2021), [homomorphic encryption] operations are about one billion times slower than normal operations."); Scheffler & Mayer, *supra* note 5, at 427.

³²¹ *See, e.g.,* Song Bian, Masayuki Hiromoto & Takashi Sato, *Towards Practical Homomorphic Email Filtering: A Hardware-Accelerated Secure Naïve Bayesian Filter*, 24 PROC. ASIA & S. PAC. DESIGN AUTOMATION CONF. 621 (2019).

³²² *See* Scheffler & Mayer, *supra* note 5, at 427–28; Théo Ryffel, David Pointcheval, Francis Bach, Edouard Dufour-Sans & Romain Gay, *Partially Encrypted Deep Learning using Functional Encryption*, 32 PROC. CONF. ON NEURAL INFO. PROCESSING SYS. 4517 (2019).

³²³ *See* Wiretap Act, 18 U.S.C. § 2511(1)(a).

³²⁴ *See* Boyden, *supra* note 191, at 702–03.

and the courts found the interception element satisfied.³²⁵ This suggests that automated content scanning at least potentially qualifies as interception under the statute.

Turning to the definition of “contents,” as discussed with respect to message franking, there is at least a plausible case that an encrypted message qualifies insofar as it is “information concerning the substance, purport, or meaning” of the message sender’s communication.³²⁶ On the one hand, the platform ideally gains no information from the encrypted material pre- or post-modification, suggesting that the encrypted message is not contents.³²⁷ On the other hand, the platform is able to manipulate and change the message’s contents, potentially even removing information from the message. These abilities can make a court more inclined to treat the encrypted message as contents rather than as metadata. If the platform can change the “substance, purport, or meaning” of a communication by manipulating homomorphically encrypted ciphertexts, then it seems reasonable to say that the ciphertexts are “information concerning” content.³²⁸

Assuming that the platform’s receipt of the message constitutes interception of contents, the platform avoids liability under the Wiretap Act only if it meets one of the statute’s exceptions: (1) the platform is the intended recipient, (2) users have consented to the interception, or (3) the automated scanning is in the platform’s ordinary course of business.

Regarding the first two exceptions, the analysis largely tracks that for message franking,³²⁹ except that the argument in favor of the exceptions is potentially weaker. One common understanding of E2EE is that the platform cannot manipulate messages based on their content—something that the platform-based content scanning techniques in fact do. If users intend the platform to intercept their messages for modification or consent to the platform modifying their messages, then that intent or consent is arguably in tension with the users’ reasonable expectations of how E2EE is supposed to work. Before concluding that users have consented, a court would likely engage in an especially searching scrutiny of a platform’s terms of service given this tension.

³²⁵ See *In re Google Inc. Gmail Litig.*, No. 13-md-2430, slip op. at 20 (N.D. Cal. Sept. 26, 2013); *In re Yahoo Mail Litig.*, 7 F. Supp. 3d 1016, 1027–28 (N.D. Cal. 2014).

³²⁶ See 18 U.S.C. § 2510(8); *supra* text accompanying notes 181–87.

³²⁷ See *DirecTV, Inc. v. Barnes*, 302 F. Supp. 2d 774, 779 (W.D. Mich. 2004).

³²⁸ See 18 U.S.C. § 2510(8).

³²⁹ See *supra* text accompanying notes 188–95.

By contrast, platform-side automated content scanning presents a stronger case for the Wiretap Act's ordinary course of business exception. Unlike message franking, where the platform sometimes retains some of the message's content (the franking or verification tags),³³⁰ platform-side scanning does not require the platform to retain any part of the encrypted message; indeed, the platform would have little reason to do so because the modified but still homomorphically encrypted message is wholly meaningless to the platform. Platform-side scanning thus appears akin to the spam detection or antivirus scanning practices that courts generally agree do not violate the statute.³³¹

3. SCA Analysis

Regarding the SCA, the only relevant point in time when communication information is stored is when the platform is processing the homomorphically encrypted content to augment or modify it. This is not a violation of § 2701 because the platform is the provider of the electronic communications service and so can authorize the processing.³³² Nor is it a violation of § 2702, because the result of the processed message is disclosed only to the intended recipient.³³³ Server-side automated content scanning thus likely avoids liability under the SCA.

4. PR/TT Analysis

For server-side scanning to be considered a PR/TT device under the PRA, it would be necessary that the scanned, homomorphically encrypted data (1) not be content and (2) include information identifying the sender or recipient of the message.³³⁴ This seems unlikely, given the above discussion of homomorphically encrypted data as contents under the Wiretap Act.³³⁵ If server-side scanning does fall within the statute, then it likely satisfies the abuse-protection exceptions, given that the purpose of scanning is to flag or otherwise affect content that the platform deems abusive.³³⁶ It likely also satisfies the service-operation and user-consent

³³⁰ See *supra* text accompanying notes 208–09.

³³¹ See *In re Google Inc. Gmail Litig.*, No. 13-2430, slip op. at 20 n.4 (N.D. Cal. Sept. 26, 2013).

³³² See Stored Communications Act, 18 U.S.C. § 2701(c)(1).

³³³ See *id.* § 2702(b)(1).

³³⁴ See Pen Register Act, 18 U.S.C. § 3127(3)–(4).

³³⁵ See *supra* text accompanying notes 326–28.

³³⁶ See 18 U.S.C. § 3121(b)(1)–(2).

exceptions.³³⁷

5. CALEA Analysis

Server-side automated content scanning likely does not implicate CALEA for at least two reasons. First, the statute regulates only voice-like telecommunications services, and absent substantial technological progress, it is unlikely that server-side scanning techniques will be applicable to real-time voice communications in the near future.³³⁸ Second, it is not clear what law enforcement would get out of asserting CALEA against platforms using server-side scanning. The statute only requires platforms to build in capabilities for intercepting information from communications,³³⁹ but the intercepted information would be encrypted and CALEA does not require the platform to decrypt it.³⁴⁰ As a result, CALEA would probably not be asserted to require modifications to the design of a server-side automated content scanning system.

6. CFAA Analysis

As a reminder, a violation under the CFAA requires intentional unauthorized access to a protected computer that results in one of several types of harm.³⁴¹ Server-side scanning does not cause most of the types of harm enumerated in the statute. The scanning platform does not “obtain[] information” because computations on homomorphically encrypted data do not reveal information to the platform.³⁴² The scanning is presumably not done “with intent to

³³⁷ See *supra* Part III.A.4.

³³⁸ More specifically, the major current limitation of homomorphically encrypted content is that computations on such content are slow, especially when the computations are complex. See *supra* note 320. Since automated analysis of real-time voice communications for content moderation would likely involve multilayer machine learning models, the computational cost of such models on homomorphically encrypted content would probably render this form of content moderation infeasible.

³³⁹ See Communications Assistance for Law Enforcement Act (CALEA) § 103(a)(1), 47 U.S.C. § 1002.

³⁴⁰ See CALEA § 103(b)(3).

³⁴¹ See Computer Fraud and Abuse Act, 18 U.S.C. § 1030.

³⁴² See *id.* § 1030(a)(2)(C). One might argue that, being derived from message information, the encrypted data is information in the same way that a cryptographic hash may be “contents” under the Wiretap Act. Cf. *supra* text accompanying notes 181–87. However, the cryptographic hash has informational value to the platform in that it can be used to provably tie a message to its sender; encrypted data is ideally indistinguishable from randomness and should not serve this

defraud” and the platform does not thereby “obtain[] anything of value” from it.³⁴³

The only plausible cause of action under the CFAA would be under § 1030(a)(5)(A), which prohibits the platform from “knowingly caus[ing] transmission of a program, information, code, or command, and as a result of such conduct, intentionally caus[ing] damage without authorization, to a protected computer.”³⁴⁴ The three candidate computers are the sender’s device, the recipient’s device, and the messaging system under a network-trespass theory. The program transmitted in all cases is the platform’s software that performs the homomorphic encryption. That software enables the platform to perform server-side scanning, which could cause “damage” under the CFAA if the scanning results in modification of the message rather than mere addition of flagging information.³⁴⁵

Since this “damage” only occurs on the recipient’s device upon receipt of a message, there is no violation with respect to the sender’s device alone. With respect to the recipient’s device standing alone, it would be hard to argue that the software on the *recipient’s device* resulted in any “damage” because it was the software on the *client’s device* that encrypted the message with a homomorphic scheme to enable server-side scanning.³⁴⁶ As a result, a violation may only occur under a network-trespass theory in which the transmissions of software to *both devices* count for purposes of the statute.

Even under the network-trespass theory, there is probably no violation of § 1030(a)(5)(A) because the damage is likely not “without authorization.” Assuming that users voluntarily installed the messaging software, the software’s actions at most exceed authorized access, which this provision of the CFAA does not prohibit.³⁴⁷ Furthermore, absent explicit consent, the messaging

informational purpose. See ROSULEK, *supra* note 24, at 22. In any event, the platform presumably has authorization to obtain the *encrypted* data insofar as the message sender desires to have the platform relay that data to the intended recipient.

³⁴³ See 18 U.S.C. § 1030(a)(4).

³⁴⁴ *Id.* § 1030(a)(5)(A).

³⁴⁵ See *id.* § 1030(e)(8) (defining damage as “impairment to the integrity or availability of data, a program, a system, or information”).

³⁴⁶ This assumes that the clause in § 1030(a)(5)(A) “knowingly causes the transmission” is modified by the trailing clause “to a protected computer.” It is not clear that this is the case; the text could be interpreted to mean that causing a transmission *anywhere* that results in unauthorized damage to a protected computer is a violation. If this is the correct interpretation of the statute, then the outcome is essentially the same as the outcome under the network trespass theory described here.

³⁴⁷ See *Int’l Airport Ctrs., LLC v. Citrin*, 440 F.3d 418, 420 (7th Cir.

parties' choice of a platform using homomorphic encryption is arguably implicit authorization for the platform to use the capabilities of homomorphic encryption. Furthermore, even in non-homomorphic encryption systems, the messaging platform can block messages from being sent based on unencrypted metadata. As a result, end-to-end encryption users should expect that platforms have some capabilities to interfere with or modify content being transmitted, mitigating the possible argument that the choice to use end-to-end encryption implies a lack of authorization for the platform to modify content.

7. POCA Analysis

Under POCA, server-side scanning does not trigger the § 2258A reporting requirements because the homomorphic encryption operations of the server do not reveal information about the contents of communications, thereby preventing any platform operator from gaining actual knowledge of reportable activity.³⁴⁸ Similarly, the immunity available under § 2258B is inapplicable because client-side scanning involves no reporting.

In the perhaps-contrived situation that a platform incorporates NCMEC's hash database into its server-side scanning system, § 2258C might come into play. Under that statute, an online service provider is permitted to use NCMEC's hash database only "for the sole and exclusive purpose of permitting that provider to stop the online sexual exploitation of children."³⁴⁹ A platform implementing server-side scanning is incapable of performing the canonical response of reporting distribution of CSAM; its use of the database must be for some other purpose such as issuing warning messages to users or blurring out contraband images.³⁵⁰ Whether such activities constitute "stop[ping] the online sexual exploitation of children" is not clear, given that this statutory language has received no judicial interpretation. It is perhaps informative, though, that POCA originally permitted providers to use the hash database only "to stop further transmission of images."³⁵¹ The current statutory language

2006).

³⁴⁸ See PROTECT Our Children Act of 2008, 18 U.S.C. § 2258A(a)(1)(A).

³⁴⁹ *Id.* § 2258C(a)(1).

³⁵⁰ If the server-side scanning of images for CSAM is coupled with client-device software that performs automatic reporting, then the system as a whole falls under the analysis of client-side scanning, as explained above.

³⁵¹ See PROTECT Our Children Act of 2008 (POCA), Pub. L. No. 110-401, § 501(a), § 2258C(a)(1), 122 Stat. 4229, 4249.

“to stop the online sexual exploitation of children” comes from amendments in the CyberTipline Modernization Act,³⁵² which perhaps suggests congressional intent to permit a broader range of platform activities using the database.

D. Client-Side Automated Content Scanning

Even in an E2EE system, the unencrypted plaintext of a message is available on the client devices that senders and recipients use. As a result, the devices themselves can scan messages for illicit content without violating users’ confidentiality expectations. Client-side scanning technologies have attracted significant attention recently, with Apple proposing one possible system and the U.K. Home Office funding the development of others.³⁵³ Indeed, several commentators have already made initial attempts at analyzing the legality of client-side scanning.³⁵⁴

The “scanning” part of client-side scanning can vary widely, from simply searching for and flagging exact matches of improper content (e.g., a wordlist-based profanity filter) to using complex perceptual hashing techniques (e.g., the PhotoDNA database for detecting child sexual abuse material).³⁵⁵ One might even envision using machine learning to build a client-side automatic content scanner that does extensive image recognition and classification. For our purposes, though, the legal analysis turns less on the particulars of the scanning algorithm and more on the series of communications

³⁵² See CyberTipline Modernization Act of 2018, Pub. L. No. 115-395, § 4(2)(A)(v), 132 Stat. 5287, 5292.

³⁵³ See Priti Patel, *I Call on the World’s Tech Giants, Please Don’t Put Profit Before Safety*, TELEGRAPH (Sept. 8, 2021, 6:00 AM), <https://www.telegraph.co.uk/politics/2021/09/08/priti-patel-call-worlds-tech-giants-please-dont-put-profit-safety/> [https://perma.cc/JR5T-SJNQ] (U.K. Home Office announcing funding for development of client-side scanning technologies for end-to-end encrypted platforms); Press Release, Dep’t for Digit., Culture, Media & Sport, Home Off., Chris Philp, Member of Parliament, Priti Patel, Member of Parliament, *Government Funds New Tech in the Fight Against Online Child Abuse* (Nov. 17, 2021), <https://www.gov.uk/government/news/government-funds-new-tech-in-the-fight-against-online-child-abuse> [https://perma.cc/YL8N-D9MF].

³⁵⁴ See, e.g., Mark Rasch, *Is Apple’s Client-Side Child Porn Scanning Legal?*, SEC. BOULEVARD (Aug. 20, 2021), <https://securityboulevard.com/2021/08/is-apples-client-side-child-porn-scanning-legal/> [https://perma.cc/X25Q-HTLF]; Paul Rosenzweig, *The Law and Policy of Client-Side Scanning*, LAWFARE (Aug. 20, 2020, 10:56 AM), <https://www.lawfareblog.com/law-and-policy-client-side-scanning> [https://perma.cc/DWV9-EU23].

³⁵⁵ See Weigel, *supra* note 13; Gernand, *supra* note 13.

between the client device and external servers. Thus, these communication protocols are the main focus of the description below.

1. Technical Overview

In this Article, we define client-side scanning as any technology in which a client device is configured to perform automated processing related to content moderation, based on communications sent to or from the device. The simplest way to implement client-side scanning in an E2EE system is to embed the entire scanner as a program or extension that runs on a user's device.³⁵⁶ Should the scanner determine that content is problematic, it can either take actions that are solely limited to the client device (e.g., displaying a warning or blurring an image), or it can transmit information about the determination to another computer or system via a network.³⁵⁷

Pure client-side scanning systems obviously face limitations: they are constrained by the user device's computing power; they often require periodic updates to catch new problematic content; savvy clients may be able to disable or modify the systems; and risks attach to giving clients access to the scanning databases and algorithms. That said, their ability to work with unencrypted content gives these systems an edge over other content moderation techniques for encrypted systems.

More advanced client-side scanning systems involve interactions with an external server. A simple option would be for the device to hash plaintext content and send the hash to a server for comparison against a database of illicit hashes. Many different architectures can be imagined, but a proposed system by Anunay Kulshrestha and Jonathan Mayer demonstrates several important features and is used as an example here.

Kulshrestha and Mayer developed a client-server protocol for determining whether a piece of content is present in a database of content (e.g., an image is an exact match for a known CSAM image), without revealing information about the contents of the database to users. First, the client with access to the unencrypted content requests a relevant portion of the database from the server, and the server responds by homomorphically encrypting the requested database portion and returning it to the client.³⁵⁸ The client then performs a

³⁵⁶ In this Part, we assume that any client-side scanning is carried out by platform-provided software on the user's device.

³⁵⁷ See Weigel, *supra* note 13, at 217 (noting both of these options as part of Apple's Communication Safety feature).

³⁵⁸ See Kulshrestha & Mayer, *supra* note 45, at 899–900. The challenge

computation using the encrypted database portion and the unencrypted content.³⁵⁹ This computation outputs a flag indicating whether the content is in the database or not.³⁶⁰

The client then sends the encrypted flag back to the server, which decrypts it. Because the flag was produced by computation on the homomorphically encrypted database portion, the flag remains encrypted, such that the client does not learn whether the content was flagged.³⁶¹ Furthermore, the computation is constructed such that the content itself cannot be discerned from the flag.³⁶² As a result, the server can take action in response to illicit content without tipping the client off, learning about clients' permissible content, or revealing information about the database of illicit content.

is requesting the relevant portion without revealing information about the content. To do this, the client homomorphically encrypts the request for the database portion, and the server applies the request to the entire database. This results in a computed result containing only the desired database portion, but since that result is encrypted, the server does not know what portion is being returned. *See id.* at 900 (lemma 8.1).

³⁵⁹ *See id.* at 900–02.

³⁶⁰ *See id.* at 902.

³⁶¹ *See id.* at 903 (theorem 11.4).

³⁶² *See id.* (theorem 11.3).

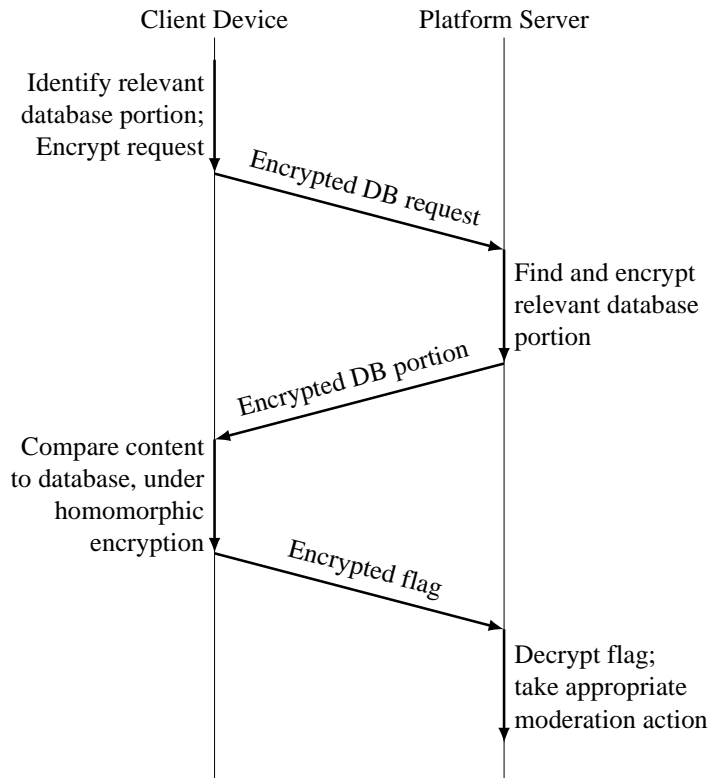


Figure 3: Diagram of communications performed in the client-side scanning system proposed by Kulshrestha and Mayer.

Other systems are essentially a subset of this process. For example, many proposed systems for detecting phishing attacks or email spam involve a client homomorphically encrypting content (e.g., a URL, email, or website screenshot) and sending it to a server that computes whether the transmitted encrypted content is undesirable. The result of that computation, still encrypted, is sent back to the client, where it can be decrypted and then acted upon without the server ever identifying its content or even the result of the computation.³⁶³ In terms of their functionality and outcomes,

³⁶³ See Edward J. Chou, Arun Guruajan, Kim Laine, Nitin Kumar Goel, Anna Bertiger & Jack W. Stokes, *Privacy-Preserving Phishing Web Page Classification via Fully Homomorphic Encryption*, IEEE INT'L CONF. ON ACOUSTICS SPEECH & SIGNAL PROCESSING 2792, 2793 fig.1 (2020) (applying homomorphic encryption to screenshots of phishing attack websites); Imtiyazuddin Shaik, Nitesh Emmadi, Harshal Tupsamudre & Harika Narumanchi, *Privacy Preserving Machine Learning for Malicious URL Detection*, 2021 DATABASE & EXPERT SYS. APPLICATIONS 31, 32 (same for URLs); Trinabh Gupta, Henrique Fingler, Lorenzo Alvisi & Michael Walfish, *Pretzel: Email Encryption and Provider-Supplied*

these systems resemble server-side scanning with homomorphic encryption.³⁶⁴ But in terms of their communications between users and platforms, they are more akin to client-side scanning, specifically the first and second transmissions of Kulshrestha and Mayer. For purposes of communication privacy law analysis, these systems more properly belong under this category.

2. Wiretap Act Analysis

Client-side automated content scanning can include two activities that are potentially relevant to the Wiretap Act. The first, which we will call a “database request,” occurs when the client device requests from the platform data to be used by the scanning algorithm, such as a relevant fragment of a matchlist database. In the second, which we call “flag computation,” the client device performs scanning computations and possibly transmits the results of those computations back to the platform. Different scanning systems may involve both, one, or neither of these activities.

The threshold question under the Wiretap Act is whether either database requests or flag computations constitute an interception. Generally, only acts contemporaneous, or at least close in time, with the overall sending of communications qualify under the statute.³⁶⁵ So, for example, if the device waits to perform the flag computation until after the message has been sent, then the flag computation would not count as “interception.”³⁶⁶ This timing distinction is notable because some client-side scanning systems, such as Apple’s Private Set Intersection protocol, require that scanning occur at the time content is received.³⁶⁷

Functions Are Compatible, PROC. CONF. ACM SPECIAL INT. GRP. ON DATA COMM’N 169, 169–70 (2017) (same for email spam detection).

³⁶⁴ See *supra* Part III.C.

³⁶⁵ See *Noel v. Hall*, 568 F.3d 743, 749 (9th Cir. 2009); *Konop v. Hawaiian Airlines, Inc.*, 302 F.3d 868, 878 (9th Cir. 2002).

³⁶⁶ The device of course must retain the message in order to perform the flag computation on it later, and that retention would qualify as interception under the Wiretap Act. But so long as the user of the device chooses to retain a copy of the message, that retention is almost certainly consented to.

³⁶⁷ See Abhishek Bhowmick, Dan Boneh, Steve Myers, Kunal Talwar & Karl Tarbe, *The Apple PSI System*, APPLE INC. 5 (2021), https://www.apple.com/child-safety/pdf/Apple_PSI_System_Security_Protocol_and_Analysis.pdf [https://perma.cc/QZD8-JZNM]. There is a good policy justification for making the timing of client-side scanning relevant to Wiretap Act liability. Where the flag computation occurs well after a message is sent or received, the user has the option of deleting the message to avoid the scanner, whereas

Assuming that the database request and/or flag computation steps qualify as interceptions, they satisfy a *prima facie* case under the Wiretap Act if the information intercepted is content. The flag computation undoubtedly involves content, because the result of the computation indicates whether the user's message is illicit under the automated scanner's content moderation standards. The database request could also be content if the request depends on information in the message, as in Kulshrestha and Mayer's proposed system. To be sure, in that system the request is homomorphically encrypted such that the platform cannot discern anything about the message from it. So, whether the database request is content will turn on the question of whether homomorphically encrypted content is content, as discussed above with respect to server-side content scanning.

To the extent that a *prima facie* Wiretap Act violation is present, we next turn to the statutory exceptions. The intended-recipient exception easily applies if there is no interaction with the platform. This exception avoids liability for automated scanning systems that, for example, only display a flag to the user or block content on the device. And where the flag computation step does inform the platform of the results, there is a strong case for the business-use exception. As discussed with respect to the extension phone cases, the business use exception can turn on whether every communication is being intercepted, as opposed to only communications involving illicit activity.³⁶⁸ Some client-side scanning systems, such as that proposed by Kulshrestha and Mayer, provide the platform with information only about flagged content, thus potentially falling within the scope of the exception as long as the grounds for flagging are tied to a legitimate business objective.

3. SCA Analysis

The relevant steps of data processing for the SCA are when (1) the user's device requests database information about illicit content, (2) the user's device scans message content, and (3) the user's device sends a report of flagged content to the server.

Section 2701. None of these steps likely gives rise to a violation of § 2701. The platform itself, as the provider of the electronic communications service over which messages are scanned, can authorize access to those communications.³⁶⁹ As with message

the user has no such choice when the computation is contemporaneous with the message transmission.

³⁶⁸ See *supra* text accompanying notes 208–09.

³⁶⁹ See Stored Communications Act, 18 U.S.C. § 2701(c)(1).

franking, the user might argue that the user's device itself is a "facility through which an electronic communication service is provided" under § 2701(a). This argument likely fails for the reasons given above.³⁷⁰

Section 2702. Steps 1 and 3 could trigger liability under § 2702 to the extent that the database request or the flagged content report goes to a third party external to the platform. Assuming that the database request contains some (possibly encrypted) content relating to the information to be scanned client-side,³⁷¹ then the question is largely the same for both points: whether divulging part of the user's content, in the process of a client-side scanning protocol, violates § 2702.

If what is being scanned is a message or communication, then the service is an electronic communication service. Assuming that the service is public, then a violation of § 2702(a)(1) occurs if the message is in "electronic storage." If the message is scanned before it is read, then the message is in "temporary, intermediate storage" that meets the statutory definition.³⁷² If the message is scanned after it is read, then under Ninth Circuit precedent, at least the message can be considered in "backup protection" storage that also meets the definition.³⁷³ However, since the unencrypted and scanned message on the user's device is probably the user's only copy of the message, there is a good argument that the message is not a backup and so, is not in electronic storage.³⁷⁴

The platform can also be a remote computing service, for example, if the user device uses the platform for cloud storage. In this case, liability under § 2702(a)(2) arises if the platform "is not authorized to access the contents" for purposes other than backup and computing. This will depend on whether the platform's terms of service provide sufficient authorization to divulge information. However, the platform's offering of end-to-end encryption is arguably in tension with the notion that the user has authorized the platform to disclose content to third parties, perhaps making a court

³⁷⁰ See *supra* notes 213–16 and accompanying text.

³⁷¹ If the request is encrypted, as in the Kulshrestha and Mayer protocol, then the database request is content to the extent that encrypted information is content as described above. See *supra* text accompanying notes 181–87. If it is not content, then there could be a violation of § 2702(a)(3) if the database is run by a government entity.

³⁷² See Wiretap Act, 18 U.S.C. § 2510(17)(A).

³⁷³ See *id.* § 2510(17)(B); *Theofel v. Farey-Jones*, 359 F.3d 1066, 1077 (9th Cir. 2004).

³⁷⁴ See *Sartori v. Schrodt*, 424 F. Supp. 3d 1121, 1133 (N.D. Fla. 2019); *Republic of the Gam. v. Facebook, Inc.*, 567 F. Supp. 3d 291, 305 (D.D.C. 2021).

more inclined to find a violation here.

Exceptions to § 2702(a) potentially do not apply. No liability occurs, for example, when the message sender or recipient (or the remote computing service subscriber) consents to content being divulged as part of a client-side scanning protocol. But as noted above, such consent is at odds with expectations of end-to-end encryption. The platform could also argue that client-side scanning is a necessary incident of the platform's service or protects the platform's rights and property.³⁷⁵ Some courts have treated this exception as giving platforms broad discretion,³⁷⁶ but by analogy to the business-use exception under the Wiretap Act, other courts may be inclined to confine this exception narrowly to necessary incidents of message transmission or storage, depending on the nature of the automated scanning.³⁷⁷ The platform may also be able to use exceptions for disclosure to law enforcement to prevent specific harms,³⁷⁸ which could justify divulging the final determination of illicit content but not divulging information in the database request at step 1 above.

4. PRA Analysis

Either the client device software or the platform server could qualify as a PR/TT device under the PRA. The determinative question is what information is transmitted during the scanning process—database requests, content hashes, or encrypted flags, for example. That information may be deemed content, exempting the client device or platform server from the statute.³⁷⁹ Even if it is not content, that information may not identify the sender or recipient of any communication, again leaving the client device or platform server outside the ambit of the statute.³⁸⁰

To the extent that a client-side scanning system does qualify as a PR/TT device, the user-consent, platform-operation, or abuse-protection exceptions of the PRA may apply.³⁸¹ Application of those exceptions will largely track the message franking analysis,³⁸² with a few notes. First, as discussed above, user consent to scanning is potentially in tension with expectations about end-to-end encryption,

³⁷⁵ See Stored Communications Act, 18 U.S.C. § 2702(b)(5).

³⁷⁶ See *Facebook*, 567 F. Supp. 3d at 309.

³⁷⁷ See *supra* text accompanying notes 203–09.

³⁷⁸ See 18 U.S.C. § 2702(b)(6)–(8).

³⁷⁹ See Pen Register Act, 18 U.S.C. § 3127(3)–(4); *supra* text accompanying notes 181–87.

³⁸⁰ See 18 U.S.C. § 3127(3)–(4).

³⁸¹ See *id.* § 3121(b).

³⁸² See *supra* Part III.A.4.

so that exception to the PRA may apply with less force.³⁸³ The abuse-protection exceptions may also depend on who possesses the device on which the automated scan occurs.³⁸⁴ Scanning content before displaying it to the recipient of a message probably presents a strong case for protecting the recipient from abuse of the service. Where a user's own content is being scanned, however, it is less plausible that client-side scanning is protecting those users from their own abuses. The platform would have to argue instead that client-side scanning protects the user community at large.³⁸⁵

5. CALEA Analysis

A client-side automated content-scanning system could be usefully modified to aid law enforcement in intercepting communications. If the system is already configured to report flagged content back to the platform, then law enforcement might demand that the platform update its matching databases or algorithms to flag certain communications of interest. For example, if law enforcement suspected a certain message platform user of money laundering, then it could ask the platform to flag the phrase "money laundering" in that user's messages as illicit content within the client-side scanning system, such that the platform would be notified of such messages and could forward that information on to law enforcement. Kulshrestha and Mayer have identified this possibility of law enforcement manipulation of client-side scanning algorithms and are developing technical strategies to act as "canaries in the coal mine," revealing whether the moderation policies have been modified.³⁸⁶

The question to be addressed here is whether law enforcement could compel such modifications under CALEA. The threshold limitation is that the statute only applies to voice-like telecommunications services. Again, automated scanning programs are unlikely to be implemented as a technical matter for real-time voice communications for a number of reasons relating to computational speed. First, these scanning programs would be

³⁸³ See Pen Register Act, 18 U.S.C. § 3121(b)(3).

³⁸⁴ See *id.* § 3121(b)(1)–(2).

³⁸⁵ The statute seems open to this interpretation: The exception provides for using a PR/TT device for the "protection of users of that service" generally. *Id.* § 3121(b)(1).

³⁸⁶ See Kulshrestha & Mayer, *supra* note 45, at 905 ("The server could also collaborate with trusted third parties (e.g., civil society groups) to validate the hash set . . ."); Sarah Scheffler, Anunay Kulshrestha & Jonathan Mayer, *Public Verification for Private Hash Matching*, PROC. IEEE SYMP. ON SEC. & PRIV. 2074 (2023).

executed on users' phones or home computers, which have limited processing power. Second, to the extent that the scanning program involves homomorphic encryption, the computational complexity likely renders real-time scanning infeasible as discussed above.

If technology for real-time scanning of voice communications does become available, then there is a good argument that law enforcement could require platforms adopting that technology to build in capabilities for flagging content of interest to law enforcement. Such capabilities would not violate CALEA's encryption exception because the computation and transmission of flagging information occurs after the content has been decrypted. To the extent that a client-side automated content scanning system works for voice communications and returns information to the platform, CALEA could be used to require the platform to modify the scanning system's content moderation policies.

The "canary" technologies that reveal whether the government has modified the content moderation policies³⁸⁷ present a further problem. Under CALEA, regulated communications services must "protect[] . . . information regarding the government's interception of communications and access to call-identifying information."³⁸⁸ So, assuming that CALEA imposes technical capability requirements on client-side scanning systems, the statute might further prohibit technologies that reveal modifications to the scanning policies.

6. CFAA Analysis

Because client-side scanning provides a platform or other entity with information from a user's computer, it potentially violates the CFAA's prohibition on unauthorizedly obtaining information from a protected computer.³⁸⁹ The user's device is the protected computer under the statute.³⁹⁰ So, the dispositive issues are (1) whether transmitting a hash or flag back to the platform constitutes obtaining information under the statute, and (2) whether the platform had authorization to access the user's computer via client-side scanning software.

Regarding the first issue, the question of whether a hash or flag constitutes "information" under the CFAA parallels the questions

³⁸⁷ See *supra* note 386.

³⁸⁸ See Communications Assistance for Law Enforcement Act (CALEA) § 103(a)(4)(B), 47 U.S.C. § 1002(a)(4)(B).

³⁸⁹ See Computer Fraud and Abuse Act, 18 U.S.C. § 1030(a)(2)(C).

³⁹⁰ See *id.* § 1030(e)(2)(B).

regarding “contents” under the Wiretap Act.³⁹¹ Unlike the latter statute, however, the CFAA offers no definition of “information.”³⁹² The few cases to have considered the statutory phrase have focused on what it means to “obtain” information, not the nature of the information itself. Nevertheless, these cases at least suggest that courts are likely to interpret “information” broadly.³⁹³ Since hash values and flags indicate something about the content on the user’s computer, they could comfortably fit within such a broad interpretation of “information.”³⁹⁴ The strongest counterargument would probably be to analogize to the Fourth Amendment, where several scholars have vigorously argued that hash values merely indicating the presence of contraband are not searches.³⁹⁵ Yet, putting aside the question of whether the CFAA reaches further than the Fourth Amendment, other scholars contend that hash matching is indeed a Fourth Amendment search, bolstering the view that hashes are information under the CFAA.³⁹⁶

Regarding the second issue, unauthorized access, the platform would point to the user’s voluntary installation of software or purchase of the device with client-side scanning software. It could also rely on its terms of service to show that the user authorized

³⁹¹ See *supra* text accompanying notes 181–87.

³⁹² See Orin S. Kerr, *Focusing the CFAA in Van Buren*, 2021 SUP. CT. REV. 155, 160 (characterizing interpretation of “obtains . . . information” as an open question under the CFAA).

³⁹³ See, e.g., *Am. Online, Inc. v. Nat’l Health Care Disc., Inc.*, 121 F. Supp. 2d 1255, 1276 (N.D. Iowa 2000) (“mere observation of the data” is sufficient for violation of 18 U.S.C. § 1030(a)(2)(C)) (quoting legislative history); *United States v. Drew*, 259 F.R.D. 449, 457 (C.D. Cal. 2009) (noting that the intentionality and obtaining-information elements of 18 U.S.C. § 1030(a)(2)(C) “will always be met when an individual using a computer contacts or communicates with an Internet website”).

³⁹⁴ While the numerical content of a hash ideally conveys no information about the underlying content, the fact that two hash values match each other strongly indicates that the underlying content is the same, which is “information” in the sense that it reduces uncertainty about the state of the world.

³⁹⁵ See, e.g., Richard P. Salgado, *Fourth Amendment Search and the Power of the Hash*, 119 HARV. L. REV. F. 38, 42 (2005); Wei Chen Lin, Comment, *Where Are Your Papers?: The Fourth Amendment, the Stored Communications Act, the Third Party Doctrine, the Cloud and Encryption*, 65 DEPAUL L. REV. 1093, 1118–19 (2016).

³⁹⁶ See, e.g., Dennis Martin, Note, *Demystifying Hash Searches*, 70 STAN. L. REV. 691, 726–27 (2018); Denae Kassotis, *The Fourth Amendment and Technological Exceptionalism After Carpenter: A Case Study on Hash-Value Matching*, 29 FORDHAM INTELL. PROP. MEDIA & ENT. L.J. 1243, 1313–14 (2019).

client-side scanning. However, a user objecting to such scanning might argue that by choosing to use end-to-end encryption, they communicated an intention to prevent others—specifically including the platform—from learning about the user’s content.³⁹⁷ The choice to use a privacy-enhancing technology, such as end-to-end encryption, would arguably be rendered a nullity if a platform could circumvent that choice with nothing more than terms of service.³⁹⁸

As one of us has argued, such questions about the scope of user consent under the CFAA are both factual and normative.³⁹⁹ From a factual perspective, the question would be whether the platform’s terms of service and other disclosures are sufficiently clear to make actual users actually understand that the client-side scanning software undercuts the end-to-end encryption as to the platform. But a court taking a more normative point of view, as some courts interpreting the CFAA have, would in effect ask whether users ought to be required to tolerate hashing or flagging of their content via client-side scanning and to incorporate that normative expectation into the meaning of “authorization” under the CFAA.⁴⁰⁰

7. POCA Analysis

The reporting requirements of POCA will likely affect client-side scanning systems only to the extent that those systems are

³⁹⁷ See Rasch, *supra* note 354 (arguing that Apple’s terms of service are insufficient to authorize client-side scanning); see also Jeffrey Vagle, *Client-Side Scanning: A New Front in the War on User Control of Technology*, JUST SEC. (Oct. 28, 2021), <https://www.justsecurity.org/78749/client-side-scanning-a-new-front-in-the-war-on-user-control-of-technology/> [<https://perma.cc/ZGA9-NL9H>] (noting implicit expectation of control over data based on ownership of the device on which it is stored).

³⁹⁸ Other commentators have observed the conflict between client-side scanning and user expectations of E2EE. See Rosenzweig, *supra* note 354 (“Aspects of the Computer Fraud and Abuse Act . . . might be read to prohibit [client-side scanning].”); Erica Portnoy, *Why Adding Client-Side Scanning Breaks End-to-End Encryption*, ELEC. FRONTIER FOUND. (Nov. 1, 2019), <https://www EFF.ORG/deep links/2019/11/why-adding-client-side-scanning-breaks-end-end-encryption> [<https://perma.cc/MK2R-CKYP>] (“Client-side scanning mechanisms will break the fundamental promise that encrypted messengers make to their users: the promise that no one but you and your intended recipients can read your messages or otherwise analyze their contents to infer what you are talking about.”).

³⁹⁹ See Grimmelmann, *Consenting*, *supra* note 123.

⁴⁰⁰ See *id.*; cf. Kerr, *Norms*, *supra* note 124 (arguing that courts should look to commonly shared norms among computer users in making such determinations).

designed to identify CSAM. If the system does not transmit information back to the platform, if it does not scan for CSAM, or if it flags CSAM and other illicit material for the platform without distinguishing the two, then the platform would lack the actual knowledge required to trigger the reporting requirement of § 2258A. If the system does identifiably flag CSAM for the platform, then the platform would have an obligation to make a report to NCMEC. But since the content of the report is “at the sole discretion of the provider,” the statute imposes no obligations on what information the client-side scanning system must transmit to the platform.⁴⁰¹ One could possibly argue that a platform’s willful blindness to CSAM gives rise to imputed actual knowledge,⁴⁰² but the explicit instruction in § 2258A that platforms need not “affirmatively search, screen, or scan” for illicit content would likely be a substantial impediment to that argument.⁴⁰³

The liability shield of § 2258B also likely has little effect on client-side scanning, since all of the activities of such scanning precede any reporting of CSAM. The immunity might be relevant, though, if the client-side scanning software itself transmits the report to NCMEC. In that case, the statute would overcome liability under the SCA’s prohibition on divulging stored communications, as described with respect to message franking.⁴⁰⁴

The database of hashes under § 2258C presents a more interesting issue with respect to client-side scanning. That statute limits platforms’ use of NCMEC’s hash database to “stop[ping] the online sexual exploitation of children.”⁴⁰⁵ If a platform implements scanning software based solely on NCMEC’s database and fulfills its reporting obligations, then the platform presumably satisfies the limitation of § 2258C. But what if the platform commingles the NCMEC database with other hashes of illicit content, and, in particular, implements a client-side scanning system like that of Kulshrestha and Mayer where the platform ultimately learns nothing beyond whether the scanned content matched something in the database?⁴⁰⁶ By limiting what information it obtains through client-side scanning, the platform may render itself unable to fulfill its obligations under § 2258C—with the caveat that those obligations are not yet defined by judicial interpretation of the statute. As a

⁴⁰¹ PROTECT Our Children Act of 2008, 18 U.S.C. § 2258A(b).

⁴⁰² *Cf. Glob.-Tech Appliances, Inc. v. SEB S.A.*, 563 U.S. 754, 769 (2011).

⁴⁰³ 18 U.S.C. § 2258A(f)(3).

⁴⁰⁴ *See supra* Part III.A.7.

⁴⁰⁵ 18 U.S.C. § 2258C(a)(1).

⁴⁰⁶ *See Kulshrestha & Mayer, supra* note 45, at 903.

result, § 2258C may require platforms to separate the NCMEC hash database from other databases or systems of content moderation, so that the NCMEC hashes can be used in a manner compliant with the statute.

IV. DISCUSSION

The above legal analysis of these new content moderation technologies offers common themes, trends, and patterns from which we can draw broader conclusions about the intersection of law and technology. Some of these broader conclusions relate to the new technologies themselves, some illuminate the specific communication privacy laws we applied, and some go to larger questions of the nature of encryption and how the law generally ought to treat it.

A. Statutory Ambiguities and Proposed Amendments

The statutory analysis in the previous Part found that the content moderation technologies in question are likely legal. But our conclusions are not drawn with strong certainty, and the previous Part is as long and detailed as it is because of the numerous statutory ambiguities and splits of authority the legal analysis must contend with.

Left unaddressed, these ambiguities could have the unfortunately ironic consequence that the communication privacy laws unintentionally reduce privacy. Online platforms have legal, ethical, and business incentives to moderate content.⁴⁰⁷ A messaging platform that hopes to moderate users' messages, then, faces a choice: either adopt end-to-end encryption and implement new technologies for content moderation, or eschew encryption and moderate content with traditional means. To the extent that those new technologies are legally risky because of interpretive ambiguities, the platform may find the latter path safer and thus, that platform's users would not enjoy the privacy benefits of encrypted messaging.

That communication privacy law is "famous (if not infamous)

⁴⁰⁷ See, e.g., James Grimmelman & Pengfei Zhang, *An Economic Model of Intermediary Liability*, 37 BERKELEY TECH. L.J. (forthcoming 2023) (manuscript at 17), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4422819 [https://perma.cc/FH4R-B7XW]; Klonick, *supra* note 1, at 1625–30; Felix T. Wu, *Collateral Censorship and the Limits of Intermediary Immunity*, 87 NOTRE DAME L. REV. 293, 301 (2011).

for its lack of clarity” is not new.⁴⁰⁸ But the ambiguities that these new content moderation technologies face go beyond run-of-the-mill interpretive questions. Instead, they reflect tensions between novel cryptographic techniques and the decades-old communications paradigms that the statutes assume.

More specifically, our analysis identified numerous elements of the communication privacy laws that are similar but not identical and that consistently pose related challenges. Almost every statute includes a consent exception, for example, but the precise rules differ across the Wiretap Act, the SCA, the PRA, and the CFAA (where it is instead called “authorization”). This Part will review several of these common statutory components to address two questions. First, can they be better interpreted in light of technological change? Second, why are they inconsistent across statutes, and would a unified modular be preferable?

1. *Information and Content*

The treatment and taxonomization of information play a role in all of the communication privacy laws. The Wiretap Act, the SCA, and the PRA draw a distinction between the “contents” of a communication and metadata about its routing.⁴⁰⁹ CALEA draws distinctions between “wire and electronic communications” on the one hand, and “call-identifying information” on the other.⁴¹⁰ The CFAA similarly proscribes “obtain[ing] information” beyond authorization, although “information” is left undefined.⁴¹¹

The content/metadata distinction does not play well with the data structures used by modern cryptographic algorithms, because it is uncertain whether “contents” under the Wiretap Act encompasses data cryptographically derived from a message, such as a hashed franking tag or homomorphically encrypted content.⁴¹² Similarly,

⁴⁰⁸ *Steve Jackson Games, Inc. v. U.S. Secret Serv.*, 36 F.3d 457, 462 (5th Cir. 1994) (citing *Forsyth v. Barr*, 19 F.3d 1527, 1542–43 (5th Cir. 1994)); see Ariana R. Levinson, *Toward a Cohesive Interpretation of the Electronic Communications Privacy Act for the Electronic Monitoring of Employees*, 114 W. VA. L. REV. 461, 463–64 (2011) (“The ECPA has been described by experts as dense, intricate, and difficult for lawmakers, lawyers, and even scholars to interpret.”) (citing sources).

⁴⁰⁹ See Wiretap Act, 18 U.S.C. § 2510(4); Stored Communications Act, 18 U.S.C. § 2702(a)(1)–(2); Pen Register Act, 18 U.S.C. § 3127(3)–(4).

⁴¹⁰ See Communications Assistance for Law Enforcement Act (CALEA) § 103(a)(1)–(2), 47 U.S.C. § 1002(a)(1)–(2).

⁴¹¹ See Computer Fraud & Abuse Act, 18 U.S.C. § 1030(a)(2)(C).

⁴¹² See *supra* text accompanying notes 181–87; *supra* text accompanying notes 326–28.

under the CFAA, there is at least a plausible question as to whether a platform “obtains . . . information” from a user’s computer when the platform receives a hash or flag indicating the illicitness of client-side scanned content.⁴¹³

These uncertainties, which commentators have noted in other legal and technological contexts,⁴¹⁴ reflect a generational divide between the statutes and the technology. In 1986, when the Electronic Communications Privacy Act amended the Wiretap Act to address electronic communications,⁴¹⁵ it would have been reasonable to assume that any content that could be usefully intercepted was readable plaintext. Encryption was known at that time,⁴¹⁶ but the expectation was that encrypted messages were “unintelligible” without the decryption keys.⁴¹⁷ There would have been little value to addressing the legal ramifications of intercepting encrypted content that had no informational use without the encryption keys.

Cryptographic hashes and homomorphic encryption, both developed primarily after the Electronic Communications Privacy Act,⁴¹⁸ disrupt this logic. The information produced by both of these

⁴¹³ See 18 U.S.C. § 1030(a)(2)(C); *supra* text accompanying notes 391–96.

⁴¹⁴ See, e.g., Paul Belonick, *Transparency is the New Privacy: Blockchain’s Challenge for the Fourth Amendment*, 23 STAN. TECH. L. REV. 114, 153 (2020) (discussing, in the context of blockchain and the Fourth Amendment, whether digital signatures are content) (citing Riana Pfefferkorn, *Everything Radiates: Does the Fourth Amendment Regulate Side-Channel Cryptanalysis?*, 49 CONN. L. REV. 1393, 1429–30 (2017)). Compare Salgado, *supra* note 395, at 42 (arguing, in the context of hard drive searches, that hash comparisons are not Fourth Amendment searches because “the hash value is no more useful than a random number”), and Lin, *supra* note 395, at 1118–19 (same), with Martin, *supra* note 396, at 726–27 (arguing that hash-based screening of emails would “approximate the use of general warrants” disallowed under the Fourth Amendment), and Kassotis, *supra* note 396, at 1313–14.

⁴¹⁵ See Electronic Communications Privacy Act of 1986, Pub. L. No. 99-508, 100 Stat. 1848.

⁴¹⁶ In 1974, the National Bureau of Standards proposed adopting a standardized algorithm for data encryption. See Encryption Algorithm for Computer Data Protection, 40 Fed. Reg. 12134 (Mar. 17, 1975).

⁴¹⁷ See ELECTRONIC COMMUNICATIONS PRIVACY ACT OF 1986, H.R. REP. NO. 99-647, at 37 (1986).

⁴¹⁸ See Bart Preneel, *The First 30 Years of Cryptographic Hash Functions and the NIST SHA-3 Competition*, TOPICS CRYPTOLOGY—CT-RSA 1, 4 (2010) (discussing development of early hash functions in the late 1980s and early 1990s); Craig Gentry, *Fully Homomorphic Encryption Using Ideal Lattices*, 41 PROC. ANN. ACM SYMP. ON THEORY

technologies is unintelligible on its own. Yet that information can be usefully intercepted and combined with other information, not to decrypt the message itself, but to authenticate the sender of a message (in the case of a digitally signed hash) or to modify or alter the underlying plaintext message (in the case of homomorphic encryption). Because the Wiretap Act did not contemplate that unintelligible encrypted content could nevertheless have informational value, the statute offers no clear answers when applied to advanced cryptographic technologies that generate informational value from encrypted content.

How could the statutory schemes better accommodate these new cryptographic techniques? One option for reform would be to deem them neither content nor metadata, thus not qualifying them for protection under any statute. This would have the benefit of preventing legal liability for these content moderation techniques because these encrypted materials could be intercepted and used without restriction. However, this is probably not an ideal result. The longer that encrypted data is retained, the more likely the underlying content will be revealed, either because advances in cryptanalysis break the encryption schemes over time or because the encryption keys fall into the hands of third parties.

Instead, it would be better for Congress to develop a statutory scheme specifically tailored to cryptographic hashes and other encrypted material. Such a statute would take into account the technical need to retain such encrypted material to facilitate content moderation while limiting its storage and distribution in view of the risks of long-term retention. The statute would thus take an intermediate approach between § 2511's strict prohibitions on the interception of content⁴¹⁹ and the SCA's and the PRA's permissiveness toward platform collection and use of metadata.⁴²⁰

2. *Consent and Authorization*

The interaction between consent and encryption provides another source of uncertainty. The Wiretap Act, the SCA, and the PRA provide exceptions based on user consent,⁴²¹ and the CFAA turns on "authorization."⁴²² Judicial decisions suggest different

COMPUT., 169, 169 (2009) (proposing the first fully homomorphic encryption scheme).

⁴¹⁹ See Wiretap Act, 18 U.S.C. § 2511(a).

⁴²⁰ See Stored Communications Act, 18 U.S.C. § 2702(a)(3); Pen Register Act, 18 U.S.C. § 3121(b)(1).

⁴²¹ See Electronic Communications Privacy Act, 18 U.S.C. §§ 2511(2)(d), 2702(b)(3), § 3121(b)(3).

⁴²² See Computer Fraud & Abuse Act, 18 U.S.C. § 1030(a)(2).

approaches to consent across the laws: courts are reluctant to infer consent under the Wiretap Act absent a strong factual showing,⁴²³ while implied consent under the CFAA appears to be found with some regularity.⁴²⁴

As discussed repeatedly above,⁴²⁵ these consent provisions highlight a fundamental tension inherent to content moderation in E2EE systems. On the one hand, a platform's terms of service can presumably authorize the platform to use a content-moderating technology, and that technology can serve important trust and safety objectives. On the other hand, a user's decision to use an E2EE platform might imply an intention to disallow the platform from accessing the user's messages or content, so that finding consent to the platform's content moderation seems to conflict with that intention.⁴²⁶

Settling (and perhaps standardizing) the requirements for consent across the communication privacy laws would of course help to reduce uncertainty and clarify the permissibility of content moderation technologies that work around end-to-end encryption. But the tension involved in the consent analyses highlights a broader question about the scope of the privacy expectations that end-to-end-encryption entails. We explore that broader question below.⁴²⁷

3. *Permitted Business Activities*

Of the laws reviewed, the PRA was the only one that specifically addressed a platform's efforts toward "protection of users . . . from abuse of service."⁴²⁸ That the communication privacy laws generally do not contemplate content moderation is unsurprising for telephone-era statutes directed to one-on-one communications. Electronic group messaging capabilities, however, create extensive opportunities for harassment, misinformation, spread of CSAM, and other forms of abuse. There is a growing perception that platforms should have an ethical duty, if not a legal one, to moderate content—

⁴²³ See *supra* text accompanying notes 188–95.

⁴²⁴ See *supra* note 239.

⁴²⁵ See *supra* text accompanying notes 188–95 (the Wiretap Act); *supra* text accompanying notes 397–400 (the CFAA).

⁴²⁶ Cf. *In re Google Inc. Cookie Placement Consumer Priv.*, 806 F.3d 125, 151 (3d Cir. 2015) (holding that users' adoption of cookie blockers "clearly communicated denial of consent"); Grimmelmann, *Spyware*, *supra* note 238, at 48–49 (noting difficulty in finding consent based on software terms of use, where software's activities conflict with other software the user has installed).

⁴²⁷ See *infra* Part IV.C.

⁴²⁸ See Pen Register Act, 18 U.S.C. § 3121(b)(1).

and platforms themselves have a compelling business reason to protect their users.

Although this Article has focused on content moderation specifically on E2EE systems, it has highlighted a general need to clarify whether and when content moderation runs afoul of the communication privacy laws. Ideally, this clarification would provide substantial room for platforms to adopt abuse-mitigation techniques. Legislatively, this could be achieved by adopting the abuse-protection exceptions from the PRA into the Wiretap Act and the SCA. But courts could achieve a similar result by clarifying that a platform's content moderation activities are within the "ordinary course of its business"⁴²⁹ and "necessarily incident to the rendition of the service."⁴³⁰

It is not clear, though, that this is the best approach for adapting the communication privacy laws to platform content moderation. Platforms often moderate content for a variety of reasons unrelated to abuse protection: promoting diversity, balancing debate viewpoints, or responding to developing emergencies, for example. From that perspective, even the exceptions in the PRA may turn out to be undesirably narrow. More importantly, if platform liability turns on whether a certain content moderation practice falls within a statutory exception, that interposes the judiciary in setting platform content moderation policies, a traditionally private matter.

4. *Computer Devices*

Although the communication privacy laws typically apply to activities on a computer or electronic device, it is often not well-defined which devices fall within their ambit, opening the door to some creative interpretations of the statutes. The Wiretap Act proscribes the use of a device to "intercept" communications, suggesting that the device should exist somewhere between the communicating parties and collect data in transit,⁴³¹ but parties have alleged, with some success, interceptions on the parties' own devices and based on the collection of data before or after message transmission.⁴³² Similarly, SCA litigants have sometimes alleged that a personal computer is a "facility" protected from unauthorized

⁴²⁹ Wiretap Act, 18 U.S.C. § 2510(5).

⁴³⁰ Stored Communications Act, 18 U.S.C. § 2702(b)(3).

⁴³¹ See 18 U.S.C. § 2510(4).

⁴³² See *In re iPhone Application Litig.*, 844 F. Supp. 2d 1040, 1062 (N.D. Cal. 2012); *In re Pharmatrak, Inc. Priv. Litig.*, 329 F.3d 9, 22 (1st Cir. 2003); *In re Google Inc. Cookie Placement Consumer Priv.*, 806 F.3d 125, 135 (3d Cir. 2015); *supra* text accompanying notes 196–202.

access under the statute.⁴³³ And the term “computer” in the CFAA, though seemingly referring to a single device,⁴³⁴ could be interpreted to encompass an entire network of computers.⁴³⁵

These creative interpretations of computer devices force courts into contorted efforts to twist the other elements of the statutory language to fit the theory of liability.⁴³⁶ But they can also give rise to unexpected pathways to liability. The network trespass theory of the CFAA, for example, was originally conceived as an argument to enhance platforms’ ability to police problematic network behavior through legal action against malicious users.⁴³⁷ Yet, the same theory potentially limits platforms’ ability to police problematic behavior through content moderation technologies, because treating an end-to-end encrypted messaging network as a single “computer” under the CFAA might make content moderation activities into unauthorized access to that “computer.”⁴³⁸ Clarifying the definition of devices across the statutes would help to avoid interpretive difficulties arising from unconventional theories of what constitutes a relevant computer.

5. *Making the Statutes More Modular*

The statutory concepts of content, authorization, computer devices, and business uses are largely common to all of the communication privacy laws. Yet, each statutory scheme introduces its own definitions and exceptions to those terms, leaving each statute with idiosyncratic and inconsistent definitions. This is perhaps most noticeable with regard to the business-use exceptions. The Wiretap Act and the SCA exempt activity that is “a necessary incident to the rendition of [a communication provider’s] service or to the protection of the rights or property of the provider of that service.”⁴³⁹ The PRA instead exempts a range of activities relating to “operation, maintenance, and testing” of a service, “protection of the rights or property of such service,” or protection “from abuse of

⁴³³ See *supra* note 214 (discussing cases).

⁴³⁴ See Computer Fraud & Abuse Act, 18 U.S.C. § 1030(e)(1) (defining “computer” as “an electronic, magnetic, optical, electrochemical, or other high speed data processing device”).

⁴³⁵ See Penney & Schneier, *supra* note 116, at 494–95.

⁴³⁶ See, e.g., *iPhone*, 844 F. Supp. 2d at 1058 (reasoning that if a user’s device is a “facility” under the SCA, then the communications provider is a “user”).

⁴³⁷ See Penney & Schneier, *supra* note 116, at 478–79.

⁴³⁸ See *supra* Part III.A.5.

⁴³⁹ Wiretap Act, 18 U.S.C. § 2510(2)(a); Stored Communications Act 18 U.S.C. § 2702(b)(5).

service or unlawful use.”⁴⁴⁰ The CFAA has no business-use exception, perhaps because it was assumed that a platform always had authorization to obtain content it handled.

Interestingly, though, most of the statutes do use a consistent definition of “contents.” This is because the SCA, PRA, and CALEA all incorporate the Wiretap Act’s definitions by reference.⁴⁴¹ Indeed, the PRA specifically states that the scope of its coverage “shall not include the contents of any communication,”⁴⁴² neatly relying on the Wiretap Act’s definition to carve up coverage between the two laws.

That model of consistency could be followed for the other statutory concepts discussed above. For example, a single definition of acceptable business uses could be incorporated into all of the communication privacy laws, simplifying interpretation and obviating the need to study each statute individually to discover one’s legal obligations.

To be sure, there may be situations where divergent definitions are desirable. The CFAA, for example, likely prohibits unauthorized obtaining of “information” rather than “contents” because the statute is intended to proscribe unauthorized metadata capture. Nevertheless, a single baseline definition of contents and metadata would still be helpful, as it would give legislators a unified set of statutory terms for defining “information” in the CFAA.

6. *CALEA Encryption Exception*

Aside from general concerns about mandated design of technical systems, CALEA presents two lines of concerns with respect to its application to the content moderation technologies discussed. First, it could result in the retention of encrypted materials for longer than would be safe.⁴⁴³ Second, the privacy guarantees of technologies like message franking depend on separation of information between the platform and messaging users,⁴⁴⁴ and technical design requirements under CALEA could vitiate that separation.

One possible way of addressing these problems is to expand CALEA’s existing exception for encryption.⁴⁴⁵ While that exception

⁴⁴⁰ See Pen Register Act, 18 U.S.C. § 3121(b).

⁴⁴¹ See Electronic Communications Privacy Act, 18 U.S.C. §§ 2711(1), 3127(1); Communications Assistance for Law Enforcement Act (CALEA) § 102(1), 47 U.S.C. § 1002.

⁴⁴² 18 U.S.C. § 3127(3)–(4).

⁴⁴³ See CALEA § 103(a)(1) (requiring platforms to enable government interception of communications “at such later time as may be acceptable to the government”).

⁴⁴⁴ See *supra* Part III.A.

⁴⁴⁵ See CALEA § 103(b)(3).

currently provides that platforms “shall not be responsible for decrypting” communications, it can further absolve platforms of requirements to intercept encrypted materials in the first place. Platforms would still retain encrypted information such as traceback records in accordance with their content moderation needs, and law enforcement would essentially enjoy the same privileges to investigate encrypted communications as the platform would enjoy to moderate those communications.

B. Insights into the Technologies

Our legal analysis of content moderation technologies also engages with a conversation about the technologies themselves. That conversation has already begun: the developers of these technologies have noted uncertainty about their own work’s normative desirability,⁴⁴⁶ commentators have debated the human rights implications of these technologies,⁴⁴⁷ and lawmakers have even introduced legislative and policy proposals on content moderation for E2EE platforms.⁴⁴⁸ But by systematically reviewing the legal elements of communication privacy and the technological elements of computer systems, we hope we have enabled a sharper focus on the specific normative questions at stake. Insofar as the

⁴⁴⁶ See, e.g., Kulshrestha & Mayer, *supra* note 45, at 905 (“We do not take a position on whether E2EE services should implement the protocols that we propose, and we have both technical and non-technical reservations ourselves.”); Tyagi et al., *supra* note 264, at 423 (“Robust policy dictating how and when to perform [message forward] tracing is necessary for protection of users’ privacy expectations.”); Issa et al., *supra* note 176, at 2337 (“[A]ny decision to use content moderation within end-to-end encrypted messengers requires weighing all of its potential benefits and risks. . . . We take no stance on the policy question . . .”).

⁴⁴⁷ See, e.g., Rosenzweig, *supra* note 354 (considering client-side scanning); BUS. FOR SOC. RESP., *supra* note 14 (Facebook-commissioned report discussing human rights implications of client-side scanning and other content moderation technologies for end-to-end encrypted platforms).

⁴⁴⁸ See Natasha Lomas, *UK Wants to Force Encrypted Platforms to Do CSAM-Scanning*, TECHCRUNCH (July 6, 2022, 8:53 AM), <https://techcrunch.com/2022/07/06/uk-osb-csam-scanning/> [<https://perma.cc/C9SM-5JLB>] (describing U.K. legislative efforts to require client-side scanning); Robert Gorwa, *European Security Officials Double down on Automated Moderation and Client-Side Scanning*, LAWFARE (June 15, 2022, 8:01 AM), <https://www.lawfareblog.com/european-security-officials-double-down-automated-moderation-and-client-side-scanning> [<https://perma.cc/3HXP-YSYZ>] (describing similar E.U. efforts).

communication privacy laws are intended to reflect federal policy on users' reasonable expectations of privacy, tensions between the law and the technologies may point to larger societal tensions.

Consider, for example, the business-use exception under the Wiretap Act.⁴⁴⁹ The case law on telephone monitoring of employee conversations⁴⁵⁰ suggests a distinction in legal treatment. Systems like server-side scanning that cryptographically manipulate content without retaining it are more likely to fall within the exception,⁴⁵¹ while systems like traceback that retain content for later use are less likely to.⁴⁵² That legal difference may reflect a policy preference for data minimization, a preference that in turn can inform the future design of content moderation technologies.

Other findings from our legal analysis similarly provide guidance for future technological development. The possibility that CALEA could enable law enforcement to gain access to stored message-franking or forward-tracing information is a useful reminder that governments can influence the content moderation process, a possibility that is already driving computer science research.⁴⁵³ Uncertainty about whether cryptographic hashes are "content" under the Wiretap Act and the SCA⁴⁵⁴ is consistent with many computer scientists' skepticism of whether hashes can be revealed to platforms without violating users' privacy expectations on end-to-end encrypted platforms.

C. What Is End-to-End Encryption?

Our analysis also raises a larger question: what even *is* end-to-end encryption in the first place? At a surface level, end-to-end encryption could be defined as a system in which a message remains encrypted all the way to its destination.⁴⁵⁵ More rigorous definitions expand on the guarantees that an end-to-end encrypted system makes, often focusing on confidentiality (i.e., unauthorized third parties cannot read messages), integrity (i.e., third parties cannot change message content), and authenticity (i.e., third parties cannot

⁴⁴⁹ See *supra* text accompanying notes 203–09.

⁴⁵⁰ See *Deal v. Spears*, 980 F.2d 1153, 1158 (8th Cir. 1992), *discussed at supra* text accompanying notes 208–09.

⁴⁵¹ See *supra* text accompanying notes 330–31.

⁴⁵² See *supra* text accompanying notes 287–91.

⁴⁵³ See, e.g., Scheffler et al., *supra* note 386.

⁴⁵⁴ See *supra* text accompanying notes 181–87.

⁴⁵⁵ See WONG, *supra* note 12, § 10.1 (defining end-to-end instant message encryption as “a concept of securing communications between two (or more) participants across an adversarial path”).

send messages purporting to be from others).⁴⁵⁶ These privacy guarantees are not made to computer systems as a technical matter, but to users, or “ends,” as a matter of system design.⁴⁵⁷

Sometimes the debate over end-to-end encryption has treated these privacy guarantees as a binary. Either they are intact or they are broken, and the starkness of the choice has made the positions of encryption advocates and critics seem wholly irreconcilable.⁴⁵⁸ Indeed, platforms’ content moderation activities on E2EE messaging systems have already spawned debates over whether the platforms “break end-to-end encryption.”⁴⁵⁹ But the technologies explored in this Article show that the privacy guarantees of E2EE can be altered in subtler ways.

Forward tracing, for example, can expose the identity of a message’s sender in limited circumstances.⁴⁶⁰ That is an alteration of the property of “deniability,” the requirement that senders of end-to-end encrypted messages cannot later be provably tied to their message content.⁴⁶¹ But it is not clear whether deniability is

⁴⁵⁶ See, e.g., Mallory Knodel, Sofia Celi, Olaf Kolkman & Gurshabad Grover, *Definition of End-to-End Encryption*, INTERNET ENG’G TASK FORCE 5, <https://datatracker.ietf.org/doc/draft-knodel-e2ee-definition/> [<https://perma.cc/CN22-QLUM>]. There are differing views as to the precise list of guarantees. See, e.g., Alec Muffett, *A Duck Test for End-to-End Secure Messaging*, INTERNET ENG’G TASK FORCE 7–8 (2021), <https://datatracker.ietf.org/doc/draft-muffett-end-to-end-secure-messaging/03/> [<https://perma.cc/C6CC-NDB5>].

⁴⁵⁷ See Britta Hale & Chelsea Komlo, *On End-to-End Encryption*, CRYPTOLOGY EPRINT ARCHIVE 6–7 (2022), <https://eprint.iacr.org/2022/449> [<https://perma.cc/5GB9-9AE2>] (analyzing the concept of “endness”); Knodel et al., *supra* note 456, at 3; Muffett, *supra* note 456, at 11.

⁴⁵⁹ See, e.g., Peter Elkind, Jack Gillum & Craig Silverman, *How Facebook Undermines Privacy Protections for Its 2 Billion WhatsApp Users*, PROPUBLICA (Sept. 7, 2021, 5:00 AM), <https://www.propublica.org/article/how-facebook-undermines-privacy-protections-for-its-2-billion-whatsapp-users> [<https://perma.cc/3NWZ-DZ8B>] (adding clarification to article on WhatsApp’s content moderation practices, to note that moderation of user-reported messages does not break encryption); Whitney Kimball, *WhatsApp Moderators Can Read Your Messages*, GIZMODO (Sept. 7, 2021), <https://gizmodo.com/whatsapp-moderators-can-read-your-messages-1847629241> [<https://perma.cc/HY3S-3QB8>] (observing “a lot of confusion about what the [Facebook] means when it says ‘end-to-end encryption’” in view of Facebook’s moderation of WhatsApp messages); Abelson et al., *supra* note 14 (challenging client-side scanning technology).
⁴⁶⁰ See *supra* Part III.B.

⁴⁶¹ See *supra* text accompanying notes 166–69.

intrinsically one of the privacy guarantees of end-to-end encryption.⁴⁶² Furthermore, forward tracing alters the deniability guarantee in only a limited way: the basic implementation only allows the platform to discover the message sender;⁴⁶³ more advanced protocols impose even stronger limits on when deniability can be overcome.⁴⁶⁴

The computer science literature has sought to taxonomize and characterize how content moderation technologies alter the privacy guarantees of end-to-end encryption. Sarah Scheffler and Jonathan Mayer, for example, distinguish between “full client privacy” technologies and “partial client privacy” ones, which offer less privacy to senders of content deemed illicit.⁴⁶⁵ Descriptively taxonomizing these alterations raises but does not answer the normative question of which alterations “break” end-to-end encryption and which are acceptable or *de minimis*. As a result, computer scientists have drawn differing conclusions as to whether particular content moderation technologies are compatible with end-to-end encryption.⁴⁶⁶

As courts, policymakers, and legal commentators assess the legality and desirability of the burgeoning range of content moderation technologies for end-to-end encrypted platforms, they will have to decide what guarantees of privacy such encryption entails—they will have to say what end-to-end encryption is. The consent provisions of the communication privacy laws offer one possible place where the law may evaluate this question, at least to the extent that a court takes into account normative considerations to

⁴⁶² See Knodel et al., *supra* note 456, at 5–6 (characterizing deniability as an “optional/desirable” feature).

⁴⁶³ This is because the platform alone has access to the full set of encrypted pointers with respect to traceback protocols, or the platform’s secret key with respect to source tracking. See Tyagi et al., *supra* note 264, at 417; Peale et al., *supra* note 272, at 1487.

⁴⁶⁴ See *supra* notes 279–82 and accompanying text.

⁴⁶⁵ See Scheffler & Mayer, *supra* note 5, at 7 tbl. 2 (characterizing literature into these categories). Scheffler and Mayer also provide other distinctions among content moderation technologies, such as server privacy and transparency. See *id.* at 4–5.

⁴⁶⁶ Compare Hale & Komlo, *supra* note 457, at 14 (“Message franking intuitively meets our definition of end-to-end encryption, because users voluntarily reveal specific messages sent to them to the service provider.”), with Scheffler & Mayer, *supra* note 5, at 415 (“Under the proposed designs of message franking the E2EE deniability property will no longer hold against the moderator . . .”). The discrepancy here arises in part because Hale and Komlo do not treat authentication as a guarantee of end-to-end encryption. See Hale & Komlo, *supra* note 457, at 12.

determine consent.⁴⁶⁷ Laws on false or deceptive advertising are another place where law may weigh in on the definition of end-to-end encryption. And generally, the ongoing debate over end-to-end encryption will need to treat such encryption not as a binary matter, but as a spectrum of privacy guarantees with subtle variations enabled by modern cryptographic content moderation technologies.

CONCLUSION

Encryption has costs, but perhaps those costs are not as severe as they have seemed to be. This, at least, is the upshot of the computer science research on content moderation in end-to-end encrypted media. We believe that legal scholars and policymakers need to take this research seriously; it reorients existing debates and opens up new lines of inquiry. And we believe that computer science researchers need to take the governing law seriously; it shapes what these systems can and cannot legally do. This Article is an exercise in taking both the “computer science” and “law” parts of “computer science and law” seriously.

On one hand, the technical details matter. End-to-end encryption is not just a black box that makes content moderation impossible. The abuse-prevention protocols we have discussed enable very specific forms of detection and reporting, and they do not fit conveniently into existing statutory boxes. Arguments over encryption regulation must engage with this detail.

On the other hand, the legal details also matter. Technologists working on encryption schemes that facilitate content moderation must navigate a surprisingly complicated legal thicket. We have seen, for example, that although the Wiretap Act, the SCA, and the PRA are broadly parallel, their statutory exceptions diverge significantly when applied to encrypted content moderation. Real-world encryption systems will have to fit within these exceptions.

At the end of the day, we are optimists. We like end-to-end encryption because we believe in privacy. We like content moderation because we believe in safety. And we are encouraged by the possibility that the world is wide enough for both of them at once—provided that the law will let them coexist.

⁴⁶⁷ See *supra* Part IV.A.2.