

# Nori: Concealing the Concealed Identifier in 5G

John Preuß Mattsson and Prajwol Kumar Nakarmi

Ericsson, Sweden

{*john.mattsson, prajwol.kumar.nakarmi*}@ericsson.com

June, 2021

## Abstract

IMSI catchers have been a long standing and serious privacy problem in pre-5G mobile networks. To tackle this, 3GPP introduced the Subscription Concealed Identifier (SUCI) and other countermeasures in 5G. In this paper, we analyze the new SUCI mechanism and discover that it provides very poor anonymity when used with the variable length Network Specific Identifiers (NSI), which are part of the 5G standard. When applied to real-world name length data, we see that SUCI only provides 1-anonymity, meaning that individual subscribers can easily be identified and tracked. We strongly recommend 3GPP and GSMA to standardize and recommend the use of a padding mechanism for SUCI before variable length identifiers get more commonly used. We further show that the padding schemes, commonly used for network traffic, are not optimal for padding of identifiers based on real names. We propose a new improved padding scheme that achieves much less message expansion for a given  $k$ -anonymity.

**Keywords**— 5G, IMSI catcher, SUPI, SUCI, IMSI, NSI, Privacy, Anonymity, Subscription Concealed Identifier, Identity Protection, Padding Scheme, Name Length Distribution

## 1 Introduction

Cellular devices such as mobile phones, tablets, and wearables have become more pervasive. Their role in the leakage of Personally Identifiable Information (PII) is also scrutinized more than ever, and rightfully so. One main category of PII consists of the permanent International Mobile Subscriber Identity (IMSI). While all pre-5G cellular networks (2G, 3G, and 4G) assigned and used temporary identifiers, the permanent IMSI was sometimes sent in cleartext over the radio interface. Active attacker could also trick cellular devices to send their IMSIs over the radio interface using so-called “IMSI catchers”, and identify as well as track victims [8, 34, 12, 19, 18]. In order to solve this problem, 5G introduced the Subscription Concealed Identifier (SUCI) mechanism used to encrypt Subscription Permanent Identifier (SUPI) as well as other protection mechanisms like strict refreshment of temporary identifiers, decoupling of the permanent identifier from the paging mechanism, and secure radio redirections [28]. SUCI is calculated by using Elliptic Curve Integrated Encryption Scheme (ECIES) [32, 33].

In this paper we discover a vulnerability (Section 4) in how 3GPP has specified the ECIES profiles for SUCI in 5G security standard TS 33.501 [7]. As SUCI encryption use AES in counter mode (CTR), the ciphertext is of same length as the input plaintext. This means means that SUCIs – even though fresh each time – could be linked to each other or with SUPIs based on their lengths. Therefore, SUCI does not provide indistinguishability for variable length identifiers since it leaks PII that can be practically useful for an attacker.

We propose padding as a fix to the above-mentioned problem. This means that SUPI is padded before encryption and the exposure of PII from SUCI is alleviated. In order to assess the effect of padding, we collected real-world name length data and used information theory metrics to empirically quantify message expansion and privacy protection. We present analysis on two sets of anonymized real-world name length data from (a) Swedish population, and (b) a multi-national Company’s employees in Sweden, China, India, and the USA. These (anonymized) datasets are representative of 5G use cases and are described in Section 3. We investigated five padding schemes inspired from [15, 27] as well as a new “tail-aware block-length” padding scheme that we designed (Section 5). For many types of distributions, the “tail-aware block-length” padding scheme is much more efficient than previous padding schemes as it in many cases provides better anonymity with less message expansion.

We call our work **Nori**<sup>1</sup>. We have presented Nori to 5G Security Task Force of GSMA’s Fraud and Security Group (FASG) [16]. GSMA FASG acknowledged that the vulnerability is real and should be fixed.

<sup>1</sup>Because 3GPP decided that SUCI is pronounced as SU-SHI.

We received positive gestures that this work will likely be pursued in 3GPP, aiming for Release-18. Also, based on the work in this paper, IETF HPKE [11] (proposed to be used in TLS 1.3 [29] and MLS [10]) has introduced a recommendation on padding.

In summary, our contributions are:

- Discover an information leakage in the SUCI construction with variable-length identifiers and propose padding as a fix.
- Propose a new tail-aware block-length padding scheme that minimizes message expansion.
- Provide real-world data on name lengths for all people in Sweden as well as a multi-national company.
- Empirical evaluation of different padding schemes on real-world name length data.
- Standardization impact: IETF HPKE has introduced recommendation on padding; 3GPP will likely introduce padding for SUCIs in Release-18.

## 2 Background

5G is the fifth generation of mobile networks standardized by 3GPP [9, 6, 2]. It is an evolution of earlier generations called 4G/LTE, 3G/UMTS, and 2G/GSM/GPRS. Each subscription in a 5G network is identified by a unique long-term identifier called Subscription Permanent Identifier (SUPI) [5]. For privacy protection of SUPI, 3GPP introduced Subscription Concealed Identifier (SUCI), specified in 3GPP TS 33.501 [7] and TS 23.003 [5].

SUCI encrypts SUPI with the Elliptic Curve Integrated Encryption Scheme (ECIES) scheme [32, 33], as shown in Fig. 1. ECIES is a hybrid scheme in which key exchange is based on asymmetric cryptography and key derivation and encryption are based on symmetric cryptography. ECIES is a probabilistic encryption scheme where the same plaintext encrypted multiple times produces completely different ciphertexts that cannot be linked to each other or the plaintext. 3GPP has standardized three protection schemes [7]: Null-scheme, Profile A, and Profile B. The “null-scheme” does not do any actual encryption, rather produces the same output as the input. The Profile A and B use Curve25519 or secp256r1 together with AES-128-CTR and HMAC-SHA-256.

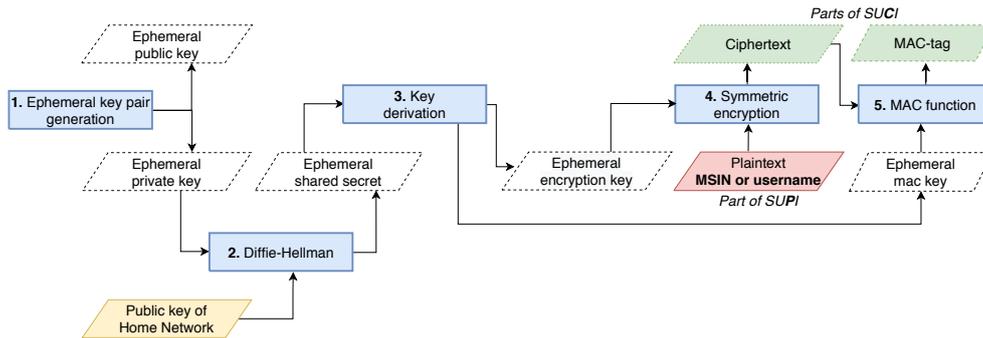


Figure 1: SUCI calculation using ECIES.

SUCI is composed as follows where the SUPI type is either IMSI or Network Specific Identifier (NSI). NSI is in the format of Network Access Identifier (NAI) as defined in IETF RFC 7542 [13]. Only the concealed identifier part is encrypted by SUCI.

$$\text{SUCI} = \text{SUPI type} \parallel \text{Home network identifier} \parallel \text{other parameters} \parallel \text{Concealed identifier}$$

When the SUPI is of type **IMSI**, the Home network identifier is composed of a 3-digit Mobile Country Code (MCC), and a 2–3-digit Mobile Network Code (MNC) and the concealed identifier contains encrypted 9–10-digit Mobile Subscription Identification Number (MSIN). When the SUPI is of type **NSI**, the Home network identifier is composed of a variable length string called the realm, and the concealed identifier contains a variable length encrypted string called the username.

## 3 Real-World Name Length Data

How NSI type SUPIs will be created in future 5G networks is yet to be seen. But it is likely that many networks will have them created from real-world names because earlier and current uses of such identifiers, e.g., in ISIMs (IP Multimedia Services Identity Module), have been based on real-world names. In this section, we present two name length data that we believe will be useful for other researchers in the field of privacy.

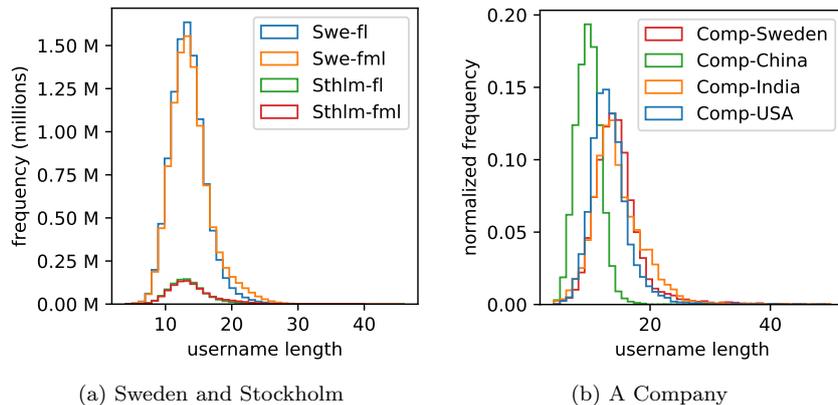


Figure 2: Histogram for name lengths. fl = first name || last name; fml = first name || maiden name || last name.

### 3.1 Name Length Data in Sweden

With the kind help from Swedish government agency SCB (Swedish: Statistiska centralbyrån) [31], we got access to the name length data for the whole of Sweden as well as Stockholm Municipality as of 31 December 2019 (Appendix A). We gathered two data sets, one with only the first and last names, and another with an additional (optional) maiden name. The name lengths are counted straight off with no spaces between different parts of the name. Their distributions, shown in Fig. 2a, are very similar, therefore, we have only analyzed the data without the maiden name in the rest of the paper. With a population of 10 million people, it is a reasonable data set to analyze the subscriber privacy of a medium-sized 5G network operator.

We note that Swedish names can contain three non-ASCII characters ‘å’, ‘ä’, and ‘ö’. These non-ASCII characters are in our analysis assumed to have a one octet encoding as they are typically replaced by ‘a’ or ‘o’ when used in NAIs and email addresses. A two octet UTF-8 encoding of ‘å’, ‘ä’, and ‘ö’ would slightly change the distributions.

### 3.2 Name Length Data in a Company

5G is to a large degree designed for use in various industries and factories. It is expected that many industries using 5G will operate their own database with subscribers and use their own realm. To that end, we gathered anonymized name length data based on email addresses of a multi-national Company. We chose four of the countries – Sweden, China, India, and the USA – it is present in, that represent different languages and cultures. Their normalized frequency distributions are shown in Fig. 2b. It can be seen that the distributions for Sweden, India, and the USA are fairly similar, while the distribution for China has much smaller mean and variance. With a few thousand employees in each country, this is a reasonable data set to analyze the subscriber privacy of a medium-sized industry using 5G.

We note that the usernames in the company email addresses had already been transliterated to ASCII characters, which is a very common practice. A UTF-8 encoding with non-ASCII character might significantly change the distribution for countries with non-roman characters, for which the recommendations in this paper might need to be adjusted.

## 4 Vulnerability in SUCI when applied to variable length SUPI

By using ECIES, SUCI achieves the notion of “indistinguishability” or “semantic security”, meaning even a very capable attacker cannot distinguish the encrypted ciphertext from a random string. However, this indistinguishability game assumes fixed length plaintexts and if the assumption is not true, then the indistinguishability is broken.

The vast majority of current 5G networks use the IMSI type SUPI where the MSIN has a fixed length for a given MCC. In that case, SUCI based on IMSI is also fixed length and provides indistinguishability. But when the NSI type SUPI is used, the username has a variable length. Therefore, SUCI based on NSI also has a variable length and indistinguishability no longer holds. An attacker gets perfect information regarding the length of the username.

Given the assumption that NSI type SUPIs will be created from real-world names, SUCI applied to the datasets in Section 3 shows poor anonymity for users with very short or very long names. In terms of  $k$ -anonymity [35], even when applied to the very large datasets from the whole of Sweden and Stockholm, the SUCI only achieves 3-anonymity. The anonymity is likely even worse as the numbers for longest length in these datasets are known to exclude data that are potentially incorrect during petition and a correct dataset

would likely only provide 1-anonymity. 1-anonymity means that an attacker can trivially identify or track at least one of the users. Similarly, SUCI applied to the company datasets from Sweden, USA, India, and China provide the worst possible 1-anonymity.

The conclusion from this is that SUCI applied to NSI type identifiers created from real-world names provides very poor anonymity for users with unusual identifier lengths, i.e., very short or very long names.

## 5 Applying Padding to Real-World Name Data

To conceal the username length leaked by SUCI and make it harder for an attacker to distinguish SUCIs based on their lengths, we suggest padding the identifier before encryption. With reference to Fig. 1, the plaintext has to be padded before Step 4.

### 5.1 Padding Schemes

We evaluated six padding schemes in total. Five of them were inspired from [27, 15]: (a) Block-length (*blk-sz-min*), (b) Power-length (*pwr-b-min*), (c) Random block-length (*rndBlk-sz-blks-min*), (d) Random-length padding (*rndLen-len*), and (e) Maximum-length (*max-len*).

The sixth padding scheme “tail-aware block-length” padding (*taBlk-l-m-r*) is designed by us. The intuition behind it is that the tails of typical distributions have the lowest frequency (meaning lower anonymity), and which benefit from padding the most. The middle parts of distributions typically have much higher frequencies and padding those only contribute to message expansion without significant increase in privacy. Therefore, we propose *taBlk-l-m-r* padding done as shown in Fig. 3, i.e., lengths below LEFT (*l*) are padded to *l*; lengths between *l* and MIDDLE (*m*) are not padded; and lengths above *m* are padded to RIGHT (*r*). By doing such selective padding, the overall message expansion is significantly reduced compared to other padding schemes that pad on all ranges of lengths.

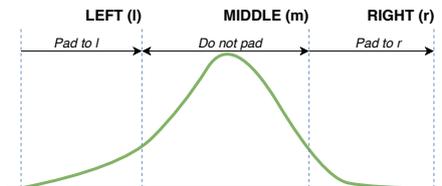


Figure 3: Illustration of the *taBlk-l-m-r* padding.

### 5.2 Evaluation of Padding Schemes

We empirically evaluated the six padding schemes on the real-world name length data described in Section 3. In total, we tested 868 padding instances with varying parameters (184 instances of *blk*, 44 of *pwr*, 184 of *rndBlk*, 32 of *rndLen*, 420 of *taBlk*, and 4 of *maxL*). For *maxL*, there is only one instance per dataset and everything is padded to the same size as the longest length in the dataset.

We assigned two attacker costs  $\alpha_1$ ,  $\alpha_2$ , and a defender cost  $\beta$  to each padding instance. The aim of the evaluation is to identify padding schemes that maximizes  $\alpha_1$  and  $\alpha_2$  while keeping  $\beta$  low. The costs are defined below:

- **Attacker Cost 1** ( $\alpha_1$ ): Conditional entropy  $H(U|P)$  [35], where  $U$  is the distribution before padding (unpadded) and  $P$  is the distribution after padding.  $\alpha_1$  represents the uncertainty about the outcome of  $U$  when the outcome of  $P$  is known, or as the expected number of bits needed to describe  $U$  given that  $P$  is known.
- **Attacker Cost 2** ( $\alpha_2$ ):  $k$ -anonymity [35], meaning the information for each person contained in a data set cannot be distinguished from at least  $k - 1$  people whose information also appear in the data set.  $k$ -anonymity is sometimes referred to as a “hiding in the crowd” guarantee.
- **Defender Cost** ( $\beta$ ): Increased bandwidth as defined in [15], i.e., a weighted sum of all the padded lengths normalized by the unpadded lengths. In this paper, we only consider the length of identifiers while [15] looks at the whole packets transported on the wire. If the identifiers are transported in larger packets, the size of padded identifiers might be small compared to the packet size.

Fig. 4 and 5 show the  $\alpha_1$ ,  $\alpha_2$  vs  $\beta$  plots for each padding instance and dataset. Each point in the plots represents a particular padding instance. Fig. 6 zooms into the plots for Sweden where  $\beta$  is limited to 2.0, i.e., maximum message size expansion of double.

### 5.3 Best Performing Padding Scheme

Just looking at Fig. 6 it clear that for the intervals shown in the figure, *taBlk-l-m-r* is performing much better than the other five padding schemes. While it is hard to put exact values on how much privacy protection

$(\alpha_1, \alpha_2)$  is needed and how much bandwidth ( $\beta$ ) is acceptable. We believe that Fig. 6 covers the values of  $\alpha_1$ ,  $\alpha_2$ , and  $\beta$  for most practical deployments. In order to choose the best padding scheme and parameters that maximize  $\alpha_1$  and  $\alpha_2$  while minimizing  $\beta$ , we use two additional metrics as defined below:

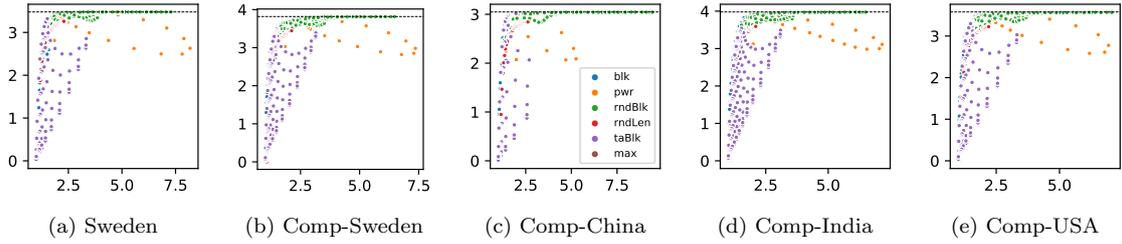


Figure 4:  $\alpha_1$  vs  $\beta$  plots.  $\alpha_1$  in  $y$ -axis.  $\beta$  in  $x$ -axis. Highest  $\alpha_1 = H(U)$  is indicated by the dashed horizontal line.

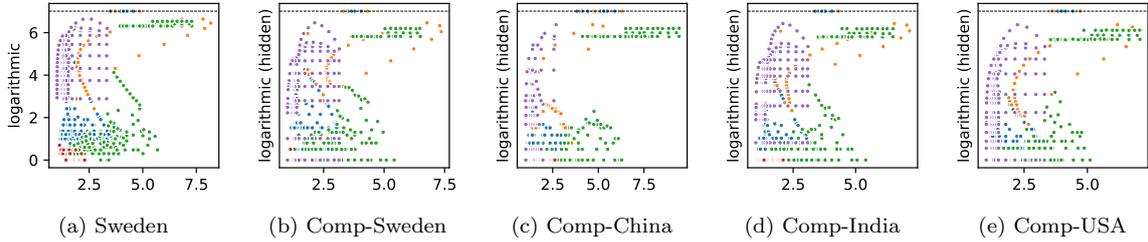


Figure 5: Log of  $\alpha_2$  vs  $\beta$  plots. Log of  $\alpha_2$  in  $y$ -axis.  $\beta$  in  $x$ -axis. Highest  $\alpha_2 = \text{population}$  is indicated by the dashed horizontal line.

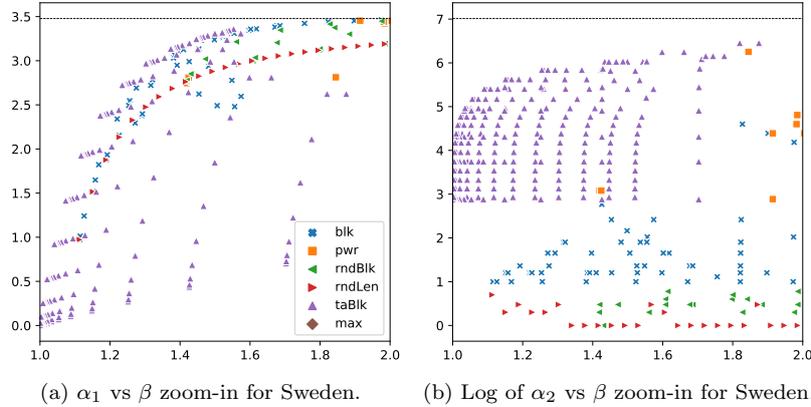


Figure 6: Zoom in plots for Sweden for  $1 \leq \beta \leq 2$ . Plot for  $max$  is not present in this range of  $\beta$ .

- **Distance to corner ( $\delta$ ):** For  $\alpha_1$  vs  $\beta$  plots in Fig. 4, the points in the top-left corner are the instances where  $\alpha_1$  is high and  $\beta$  is low. Therefore, the points with low  $\delta$  are better. Fig. 7 shows the best performing padding schemes that have lowest  $\delta$  for maximum  $\beta$  of 2. It can be seen that  $taBlk$  performs consistently better than others in terms of maintaining lower  $\delta$ . The distance metric used is the euclidean distance in the plotted 2-space, i.e., if the maximum  $\alpha_1$  value is 4.0, the point (1.0, 3.0) is assigned the distance  $\delta = \sqrt{2}$ . We note that the chosen distance metric is not necessarily the best; there could be other distance metrics giving different weights to  $\alpha_1$  and  $\beta$ , and therefore giving different results.
- **Threshold anonymity:** For  $\alpha_2$  vs  $\beta$  plots in Fig. 5, it could be sufficient to achieve some threshold anonymity  $\alpha_2$ . Fig. 8 shows an example in which the threshold  $\alpha_2$  is set to 100 and lowest  $\beta$  required by padding schemes are identified. In this case too,  $taBlk$  consistently performs better than others in terms of requiring lower  $\beta$ . We note that 100-anonymity may not provide acceptable anonymity in cases where an attacker has access to other out-of-band information.

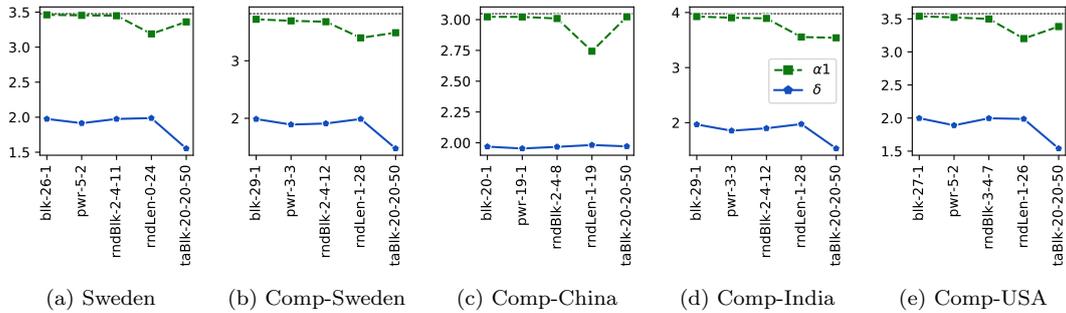


Figure 7: Best performing padding schemes in terms of lowest  $\delta$  for  $\beta \leq 2$ .

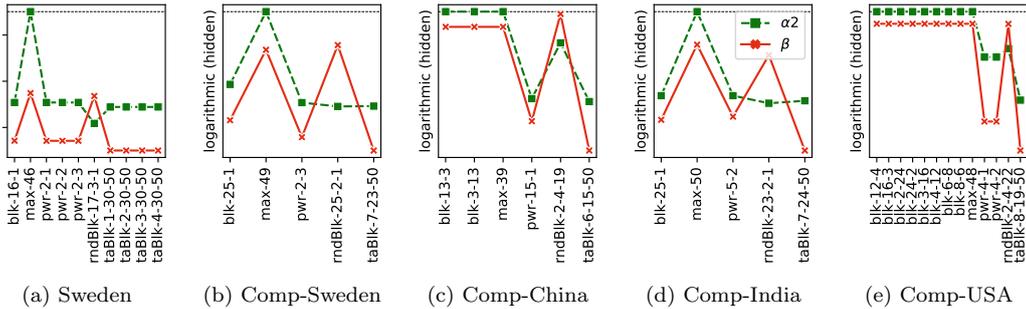


Figure 8: Best performing padding schemes in terms of lowest  $\beta$  for  $\alpha_2 \geq 100$ .

## 6 Recommendations

**Standardization.** We strongly recommend 3GPP and GSMA to mandate padding for all future SUCI generation profiles in 3GPP TS 33.501 [7]. For such future schemes, we recommend including SUPI of both types (IMSI and NSI) for padding, because although IMSI is fixed length today, it may need expansion in future. We also recommend to introduce padding enabled versions of existing Profile A and B, because it is possible to do so in a backward compatible way. We further recommend that 3GPP and GSMA specify requirements and guidelines on how the impact of the Home network private key leakage can be mitigated, e.g., by limiting the number of subscribers that use a single public key (corresponding to the private key), and changing the private keys frequently.

**Padding Scheme.** We recommend the tail-aware block-length padding (*taBlk-l-m-r*) scheme because of its superior performance in terms of  $\alpha_1$ ,  $\alpha_2$ ,  $\beta$ , and  $\delta$ . The choice of parameters ( $l, m, r$ ) depends on name length distribution and how much message expansion ( $\beta$ ) can be tolerated. *taBlk-l-m-r* instances with  $\beta \approx 2.0$  give good performance and excellent privacy. Furthermore, for Swedish name length data (Fig. 6b), *taBlk-l-m-r* instances with  $\beta$  just above 1.0 (almost no average message expansion) gives 10000-anonymity which might be good enough for some deployments. It is hard to say what a minimum acceptable  $k$ -anonymity for 3GPP deployments are as an attacker might use additional out-of-band information like realm, location, time, and cell tower identifier in identifying subscribers. A realm that is used over a large geographical area should have a much higher  $k$ -anonymity than a realm used in a small geographical area. As it is difficult to analyze all the out-of-band information that an attacker might have access to, we recommend that the targeted level of  $k$ -anonymity is chosen with large margins. For roaming users and small companies with their own realm, good anonymity is likely not possible even with padding. If non-ASCII characters are used in NAIs without transliteration to ASCII, the distributions given in this paper might change substantially. In particular such distribution could in the worst case have low (non-zero) frequencies in the middle of the distribution. A slightly modified tail-aware block-length padding with a bit block size larger ( $> 1$ ) in the middle would likely be an optimal choice also for those distributions.

According to the type of radio access technology, the messages carrying SUCI have maximum size from 1600 to 9000 bytes [3, 4, 1]. In all our datasets, the maximum username length we observed is 50. Given this and other sources (e.g., maximum 70 in e-GIF [14] and 64 in SMTP [23]), accepting  $\beta \approx 2.0$  seems to be safe and still allow plenty of space for padding.

## 7 Ethics and Responsible Disclosure

The data from Sweden and Stockholm were already anonymized by SCB; the authors only received a frequency table stating how many people have a name with a certain length. Similarly, the data from the multinational company were anonymized to contain only the lengths of the username part of the email addresses.

To anonymize the Company’s data further, the name of the Company is not revealed in this paper, and corresponding plots are either normalized or shown in a logarithmic scale.

We have done responsible disclosure of our vulnerability discovery and solution proposal to the 5G Security Task Force of GSMA’s Fraud and Security Group (FASG) [16]. GSMA FASG acknowledged that the vulnerability is real and should be fixed. We received positive gestures that this work may be pursued in 3GPP, aiming for Release-18.

## 8 Related Work

Real-world data sets of name lengths are surprisingly hard to find in the open literature. Only a few earlier works [17, 20] have investigated the statistical distribution of the length of names.

IPsec ESP packets [22] and TLS 1.3 records [29] provide mechanisms to add padding to the plaintext before encryption. They however do not provide any padding policies, leaving that to the implementation. TLS Encrypted Client Hello [30] recommends a padding scheme where small identities are padded to a fixed length and long identities are padded to the nearest multiple of 32 bytes. Padding for DNS request and response messages are specified in [26] based on empirical research and recommendations in [15, 27]. In [21], the authors discussed using padding schemes – OAEP for RSA and ISO 10126 for AES. The authors did not suggest any padding for ECIES since the encryption key itself is different every time, therefore, the purpose of padding in [21] is to ensure the so-called probabilistic encryption and does not cater privacy leakage through length of encrypted identifier.

In [24], the authors propose a formal method for quantifying the amount of privacy protection provided by traffic padding solutions. The model encompasses the privacy requirements, padding costs, and padding methods. This enables applications to choose more optimal padding methods, minimizing padding costs for a given amount of privacy, or maximizing the amount of privacy for a given padding cost. Wagner et al. [35] give an overview of privacy metrics that have been proposed in the literature and classify them based on the aspect of privacy they measure. They conclude that the lack of standardized privacy metrics makes it hard makes the choice of privacy metric and comparison between different studies hard.

Mathur et al. [25] analyse the privacy protection given by padding and encryption of traffic data sent over a network. They conclude that previous work has been too focused on padding to a fix length, which waste unnecessary amounts of bandwidth. They recommend a randomized padding and use conditional entropy as a privacy metric to evaluate the protection.

## 9 Conclusion and Further Work

5G SUCI uses state-of-the-art cryptography and gives very good privacy protection when used on fixed length identifiers such as IMSI. But when used on variable length identifiers such as NSI, it leaks significant amount of information that may be practically useful for an attacker to identify and track a victim. We presented real-world name length data, discussed padding as a solution, and showed empirical evaluation of several padding schemes. We have presented our work to GSMA FSAG’s 5G Security Task Force [16]. We strongly recommend 3GPP and GSMA to pursue this work and standardize a padding mechanism for all SUPI types in 5G.

Future research should consider the name distributions in which usernames are not transliterated to ASCII characters; such distributions could in the worst case have low frequencies in the middle of the distribution and require a slightly different padding scheme. It might also prove important to measure the message sizes that the identifiers are transported in; if these messages are large, the bandwidth overhead as measured in percentage of the original message could be small and more padding might be acceptable. Future investigations are also necessary on how large typical 5G realms will be, how much additional information (e.g., realm, location, time, and other data) attackers can use, and therefore what acceptable minimum values for  $k$ -anonymity in various deployments should be.

## Acknowledgements

We thank Ann-Marie Persson for enabling our research with SCB data and Daniel Kahn Gillmor for providing technical comments. We also thank our Ericsson colleagues for reviewing the manuscript and especially Erik Thormarker for brainstorming ideas with us.

## References

- [1] 3GPP. Evolved Universal Terrestrial Radio Access (E-UTRA); Packet Data Convergence Protocol (PDCP) specification. TR 36.323.
- [2] 3GPP. NG-RAN; Architecture description. TS 38.401.
- [3] 3GPP. Non-Access-Stratum (NAS) protocol for 5G System (5GS); Stage 3. TR 24.501.
- [4] 3GPP. NR; Packet Data Convergence Protocol (PDCP) specification. TR 38.323.

- [5] 3GPP. Numbering, addressing and identification. TS 23.003.
- [6] 3GPP. Procedures for the 5G System (5GS). TS 23.502.
- [7] 3GPP. Security architecture and procedures for 5G System. TS 33.501.
- [8] 3GPP. Study on the security aspects of the next generation system. TR 33.899.
- [9] 3GPP. System architecture for the 5G System (5GS). TS 23.501.
- [10] BARNES, R., BEURDOUCHE, B., MILLICAN, J., OMARA, E., COHN-GORDON, K., AND ROBERT, R. The Messaging Layer Security (MLS) Protocol. Internet-Draft draft-ietf-mls-protocol-11, Internet Engineering Task Force, Dec. 2020. Work in Progress.
- [11] BARNES, R., BHARGAVAN, K., LIPP, B., AND WOOD, C. A. Hybrid Public Key Encryption. Internet-Draft draft-irtf-cfrg-hpke-08, Internet Engineering Task Force, Feb. 2021. Work in Progress.
- [12] BORGAONKAR, R., HIRSCHI, L., PARK, S., AND SHAIK, A. New Privacy Threat on 3G, 4G, and Upcoming 5G AKA Protocols. *Proceedings on Privacy Enhancing Technologies 2019* (2018), 108 – 127.
- [13] DEKOK, A. The Network Access Identifier. RFC 7542, May 2015.
- [14] E-GOVERNMENT INTEROPERABILITY FRAMEWORK (E-GIF). Volume 2 – Data Types Standards, 2015.
- [15] GILLMOR, D. K. Empirical DNS Padding Policy, 2017.
- [16] GSMA. GSMA Security, 2021.
- [17] HEALY, M. The Lengths of Surnames. *Journal of the Royal Statistical Society: Series A (General)* 131, 4 (1968), 567–568.
- [18] HUSSAIN, S. R., ECHEVERRIA, M., CHOWDHURY, O., LI, N., AND BERTINO, E. Privacy Attacks to the 4G and 5G Cellular Paging Protocols Using Side Channel Information. In *NDSS* (2019).
- [19] HUSSAIN, S. R., ECHEVERRIA, M., KARIM, I., CHOWDHURY, O., AND BERTINO, E. 5GReasoner: A Property-Directed Security and Privacy Analysis Framework for 5G Cellular Network Protocol. *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security* (2019).
- [20] JACKSON, O. A., BOYETEEY, D. B., BEMILE, R. K., AND ACHEAMPONG, E. Fitting a Distribution for Length of Last names in a mixed African community in Ghana. *Open Science Repository Mathematics*, open-access (2013), e23050440.
- [21] KAROPOULOS, G., KAMBOURAKIS, G., GRITZALIS, S., AND KONSTANTINOU, E. A framework for identity privacy in sip. *J. Netw. Comput. Appl.* 33 (2010), 16–28.
- [22] KENT, S. IP Encapsulating Security Payload (ESP). RFC 4303, Dec. 2005.
- [23] KLENSIN, D. J. C. Simple Mail Transfer Protocol. RFC 5321, Oct. 2008.
- [24] LIU, W. M., WANG, L., CHENG, P., REN, K., ZHU, S., AND DEBBABI, M. Pptp: Privacy-preserving traffic padding in web-based applications. *IEEE Trans. Dependable Secur. Comput.* 11 (2014), 538–552.
- [25] MATHUR, S., AND TRAPPE, W. Bit-traps: Building information-theoretic traffic privacy into packet streams. *IEEE Transactions on Information Forensics and Security* 6, 3 (2011), 752–762.
- [26] MAYRHOFER, A. The EDNS(0) Padding Option. RFC 7830, May 2016.
- [27] MAYRHOFER, A. Padding Policies for Extension Mechanisms for DNS (EDNS(0)). RFC 8467, Oct. 2018.
- [28] NAKARMI, P. K., HENDA, N. B., AND TSIATSI, V. 3GPP Release 15: An end to the battle against false base stations?, 2019.
- [29] RESCORLA, E. The Transport Layer Security (TLS) Protocol Version 1.3. RFC 8446, Aug. 2018.
- [30] RESCORLA, E., OKU, K., SULLIVAN, N., AND WOOD, C. A. TLS Encrypted Client Hello. Internet-Draft draft-ietf-tls-esni-09, Internet Engineering Task Force, Dec. 2020. Work in Progress.
- [31] SCB. Statistics Sweden., 2021.
- [32] SECG. SEC 1: Recommended Elliptic Curve Cryptography.
- [33] SECG. SEC 2: Recommended Elliptic Curve Domain Parameters.
- [34] SHAIK, A., SEIFERT, J.-P., BORGAONKAR, R., ASOKAN, N., AND NIEMI, V. Practical Attacks Against Privacy and Availability in 4G/LTE Mobile Communication Systems. *ArXiv abs/1510.07563* (2016).
- [35] WAGNER, I., AND ECKHOFF, D. Technical privacy metrics: A systematic survey. *ACM Comput. Surv.* 51, 3 (June 2018).

## A SCB Data

Table 1: Frequency tables for name lengths in Sweden and Stockholm – as of 31 December 2019 – in the forms first name || last name (fl) and first name || maiden name || last name (fml). Source: Swedish government agency SCB (Statistiska centralbyrån) [31].

| Length | Swe-fl  | Swe-fml | Sthlm-fl | Sthlm-fml |
|--------|---------|---------|----------|-----------|
| 4      | 770     | 753     | 98       | 97        |
| 5      | 4815    | 4709    | 574      | 562       |
| 6      | 15053   | 14380   | 1717     | 1621      |
| 7      | 61086   | 57613   | 6375     | 5897      |
| 8      | 198946  | 187307  | 20511    | 18821     |
| 9      | 466161  | 440204  | 47998    | 44361     |
| 10     | 845220  | 800046  | 85583    | 79169     |
| 11     | 1232128 | 1168325 | 119605   | 110603    |
| 12     | 1537415 | 1459671 | 142089   | 131637    |
| 13     | 1634987 | 1554368 | 145003   | 134628    |
| 14     | 1441337 | 1375360 | 124717   | 116632    |
| 15     | 1072834 | 1034405 | 92960    | 88677     |
| 16     | 696269  | 688363  | 62391    | 62102     |
| 17     | 426131  | 446300  | 40174    | 43485     |
| 18     | 256513  | 300180  | 25909    | 32212     |
| 19     | 151167  | 210269  | 17260    | 25176     |
| 20     | 92937   | 157974  | 12102    | 20763     |
| 21     | 59309   | 122712  | 8363     | 16316     |
| 22     | 38522   | 92442   | 5881     | 12548     |
| 23     | 24537   | 66665   | 3961     | 9082      |
| 24     | 15231   | 44551   | 2459     | 5978      |
| 25     | 9243    | 28065   | 1471     | 3725      |
| 26     | 5650    | 16943   | 974      | 2285      |
| 27     | 3497    | 9977    | 579      | 1351      |
| 28     | 2145    | 5775    | 371      | 801       |
| 29     | 1365    | 3373    | 235      | 480       |
| 30     | 864     | 1969    | 160      | 294       |
| 31     | 554     | 1183    | 99       | 186       |
| 32     | 378     | 719     | 72       | 127       |
| 33     | 241     | 437     | 53       | 82        |
| 34     | 186     | 301     | 40       | 58        |
| 35     | 137     | 189     | 37       | 44        |
| 36     | 74      | 119     | 11       | 20        |
| 37     | 42      | 61      | 5        | 9         |
| 38     | 41      | 53      | 11       | 15        |
| 39     | 25      | 30      | 3        | 3         |
| 40     | 21      | 25      | 22       | 26        |
| 41     | 14      | 20      | 0        | 0         |
| 42     | 14      | 14      | 0        | 0         |
| 43     | 8       | 10      | 0        | 0         |
| 44     | 7       | 8       | 0        | 0         |
| 45     | 3       | 5       | 0        | 0         |
| 46     | 13      | 17      | 0        | 0         |