

The 9 Lives of Bleichenbacher’s CAT: New Cache ATtacks on TLS Implementations

Abstract—At CRYPTO’98, Bleichenbacher published his seminal paper which described a padding oracle attack against RSA implementations that follow the PKCS #1 v1.5 standard. Over the last twenty years researchers and implementors had spent a huge amount of effort in developing and deploying numerous mitigation techniques which were supposed to plug all the possible sources of Bleichenbacher-like leakages. However, as we show in this paper most implementations are still vulnerable to several novel types of attack based on leakage from various microarchitectural side channels: Out of 9 popular implementations of TLS that we tested, we were able to break the security of 7 implementations with practical proof-of-concept attacks. We demonstrate the feasibility of using those Cache-like ATtacks (CATs) to perform a downgrade attack against any TLS connection to a vulnerable server, using a BEAST-like Man in the Browser attack. The main difficulty we face is how to perform the thousands of oracle queries required before the browser’s imposed timeout (which is 30 seconds for almost all browsers, with the exception of Firefox which can be tricked into extending this period). The attack seems to be inherently sequential (due to its use of adaptive chosen ciphertext queries), but we describe a new way to parallelize Bleichenbacher-like padding attacks by exploiting any available number of TLS servers that share the same public key certificate. With this improvement, we could demonstrate the feasibility of a downgrade attack which could recover all the 2048 bits of the RSA plaintext (including the premaster secret value, which suffices to establish a secure connection) from 5 available TLS servers in under 30 seconds. This sequential-to-parallel transformation of such attacks can be of independent interest, speeding up and facilitating other side channel attacks on RSA implementations.

I. INTRODUCTION

The Public Key Cryptography Standard #1 (PKCS #1) [53] is the main standard used for implementing the RSA public key algorithm [51] in a large variety of security protocols. Twenty years ago, Bleichenbacher [10] demonstrated that the padding scheme defined in PKCS #1 v1.5 (which is used to map shorter messages into full length RSA plaintexts) is vulnerable to a *padding oracle* attack. Specifically, given an indication whether the plaintext which corresponds to a given ciphertext is correctly formatted, an attacker can mount an adaptive chosen ciphertext attack which recovers the full plaintext from any given ciphertext.

Since its publication, multiple Bleichenbacher-like attacks have been demonstrated, exploiting a large variety of oracles, including error messages [11, 36], timing variations [34, 42] and memory access patterns [62]. After each attack, implementors had adopted some mitigation techniques to ensure that the use of PKCS #1 v1.5 does not leak information on the padding, but these mitigation techniques had become increasingly difficult to understand, to implement, and to maintain. Considering the number of demonstrated attacks and

the ongoing mitigation efforts, we set out in this paper to answer the following basic question:

Are modern implementations of PKCS #1 v1.5 secure against padding oracle attacks?

A. Our Contribution.

Regrettably, our answer to this question is negative, as the vast majority of implementations we evaluated are still vulnerable to padding oracle attacks. Making the situation worse, we show that padding oracle attacks can be made extremely efficient, via more careful analysis and novel parallelization techniques. Finally, we show that while the use of RSA key exchange is declining, padding oracles can be used to mount downgrade attacks, posing them as a threat to the security of a much larger number of connections (including those done via protocols that do not even support the RSA key exchange).

More specifically, our contributions are as follows.

New Techniques for Microarchitectural Padding Oracle Attacks. We have tested nine fully patched implementations of various RSA-based security protocols (OpenSSL, Amazon s2n, MbedTLS, Apple CoreTLS, Mozilla NSS, WolfSSL, GnuTLS, BearSSL and BoringSSL). While all of these implementations attempted to protect against microarchitectural and timing side channel attacks, we describe new side channel attack techniques which overcome the padding oracle countermeasures. Notably, out of the 9 evaluated implementations, only the last two (BearSSL and BoringSSL) could not be successfully attacked by our new techniques.

Downgrade Attacks. Next, we show the feasibility of performing downgrade attacks against all the deployed versions of TLS, including the latest TLS 1.3 standard (which does not even support RSA key exchange). More specifically, even though the use of RSA in secure connections is diminishing (only $\approx 6\%$ of TLS connections currently use RSA[1, 45]), this fraction is still too high to allow vendors to drop this mode. However, as we show in [Section VI](#), supporting this small fraction of users puts everyone at risk, as it allows the attacker to perform a downgrade attack by specifying RSA as the only public key algorithm supported by the server.

Attack Efficiency. Rather than targeting premaster secrets of individual connections, we adopt a BEAST-like[20] approach, targeting instead the long term login tokens. As only a single broken connection is sufficient to recover the login token, in [Section VI](#) we show that the query complexity of padding oracle attacks can be substantially reduced (at the expense of the success probability of breaking a specific connection), while still preserving the attacker’s ability to extract login

tokens before the connection timeout enforced by almost all web browsers.

Attack Parallelization. As another major contribution, we show in this paper a novel connection between padding oracle attacks and the Closest Vector Problem (CVP). While there are some techniques to parallelize parts of padding oracle queries [36], those techniques could not overcome the fact that (perfect) padding oracles attacks use adaptively chosen ciphertexts, thereby requiring that some parts of the attack remain sequential. Using lattice reduction techniques we overcome this limitation, by combining the results from different (and partial) attack runs against different servers sharing the same RSA key. With those techniques in hand, we were able to show for the first time the feasibility of recovering a full 2048-bit RSA plaintext using 5 TLS servers in under the 30 second timeout enforced by almost all web browsers.

B. Software Versions and Responsible Disclosure

Our attacks were performed on the most updated versions of the cryptographic libraries evaluated, as published at the time of writing. We have compiled each library using its default compilation flags, leaving all side channel countermeasures in place. Following the practice of responsible disclosure, we disclosed our findings in August 2018 to all the vendors mentioned in this paper. We have further participated in the design and empirical verification of the proposed countermeasures. Updated versions of the affected libraries will be published in a coordinated public disclosure in November 2018. We note that, coincidentally, OpenSSL partially patched one of the vulnerabilities we discovered, independently and in parallel to our disclosure process.

II. BACKGROUND

A. Padding Oracle Attacks on TLS

TLS has a long history of padding oracle attacks of different types. Those attacks led to the development and implementation of new mitigation techniques, and then new attacks.

The Lucky 13 attack by AlFardan and Paterson [5] showed how to use a padding oracle attack to break TLS CBC HMAC encryption. Irazoqui et al. [33] and Ronen et al. [52] have shown how to use cache attacks to attack code that has been patched against the original attack.

After the publication of the Bleichenbacher attack, the TLS specifications defined a new mitigation with the goal of removing the oracle [16, 17, 18]. However, it seems that completely removing the oracle is a very difficult task as was shown by multiple cycles of new attacks and new mitigations [11, 36, 42]. As we show in our paper, Bleichenbacher type attacks are still possible even on fully patched implementations.

B. RSA PKCS #1 v1.5 Padding

In this section we describe the PKCS #1 v1.5 padding standard, which dictates how a message should be padded before RSA encryption. Let (N, e) be an RSA public key, let (d, p, q) be the corresponding private key, and let ℓ be

the length of N (in bytes). The encryption of a message m containing $k \leq \ell - 11$ bytes is performed as follows.

- 1) First, a random padding string PS of byte-length $\ell - 3 - k \geq 8$ is chosen such that PS does not contain any zero-valued bytes.
- 2) Set m^* to be $0x00||0x02||PS||0x00||m$. Notice that the length of m^* is exactly ℓ bytes.
- 3) Interpret m^* as an integer $0 < m^* < N$ and compute the ciphertext $c = m^{*e} \bmod N$.

The decryption routine computes $m' = c^d \bmod N$ and parses m' as a bit string. It then checks whether m' is of the form $m' = 0x00||0x02||PS''||0x00||m''$ where PS'' is a string consisting of at least 8 bytes, all of them must be non-zero. In case this condition holds the decryption routine returns m'' . Otherwise the decryption routine fails.

C. Bleichenbacher's Attack on PKCS #1 v1.5 Padding

In this section we provide a high level description of Bleichenbacher's "million message" attack [10] on the PKCS #1 v1.5 padding standard described above. At a high level, the attack allows an attacker to compute an RSA private key operation (e.g., $m^d \bmod N$) on a message m of his choice without knowing the secret exponent d .

Attack Prerequisites. Bleichenbacher's attack assumes the existence of an oracle Bl which on input a ciphertext c answers whether c can be successfully decrypted using RSA PKCS #1 v1.5 padding as described above. More formally, let (N, d) be an RSA private key. The oracle Bl performs the following for every ciphertext c

$$Bl(c) = \begin{cases} 1 & \text{if } c^d \bmod N \text{ has a valid PKCS \#1 v1.5 padding} \\ 0 & \text{otherwise} \end{cases}$$

As was previously shown, such an oracle can be obtained by several types of side channel leakage [11, 34, 36, 42, 62]. We now proceed to describe how an attacker can utilize the Bleichenbacher oracle Bl in order to perform an RSA secret key operation on c (such as decryption or signature) without knowing the secret exponent d . We refer the reader to [10] for a more complete description.

High Level Attack Description. Let c be an integer. In order to compute $m = c^d \bmod N$, the attack proceeds as follows.

- **Phase 1: Blinding.** The attacker chooses different random integers s_0 and computes $c^* \leftarrow c \cdot s_0^e \bmod N$. The attacker checks if c^* is a valid PKCS #1 v1.5 ciphertext by evaluating $Bl(c^*)$. This phase terminates when a suitable s_0 is found such that $Bl(c^*) = 1$ and can be skipped completely in case c is already a valid PKCS #1 v1.5 ciphertext (i.e. $s_0 = 1$). We notice that as soon as the oracle succeeds $Bl(c^*) = 1$ the attacker knows that the corresponding message $m^* = m \cdot s_0 = m^{*ed} \bmod N$ starts with $0x0002$. Thus, it holds that $m \cdot s_0 \bmod N \in [2B, 3B)$ where $B = 2^{8(\ell-2)}$ and ℓ is the length of N in bytes. Finally, the condition of $m \cdot s_0 \bmod N \in [2B, 3B)$ implies that there exists an integer r such that $2B \leq m \cdot s_0 - rN < 3B$, or equivalently:

$$(2B + rN)/s_0 \leq m < (3B + rN)/s_0.$$

- **Phase 2: Range Reduction.** Having established that $\frac{2B+rN}{s_0} \leq m < \frac{3B+rN}{s_0}$, the attacker proceeds to choose a new random integer s , computes $c^* \leftarrow c \cdot s^e \bmod N$ and checks that $\text{Bl}(c^*) = 1$. When a suitable s is found, the adversary can further reduce the possible ranges of m , see [10] for additional details. The attack terminates when the possible range of m is reduced to a single candidate.

Attack Efficiency. For N consisting of 1024-bits, Bleichenbacher’s original analysis [10] requires about 1 million calls to the oracle Bl (e.g., requiring the attacker to observe 1 million decryptions). However, subsequent analysis had shown that this number can be improved to as little as 3800 oracle queries for realistic scenarios [7].

The Noisy Oracle Case. We notice here that Bleichenbacher’s attack does not require the oracle Bl to be perfect and is capable of handling one sided errors where $\text{Bl}(c) = 0$ while c is a valid PKCS #1 v1.5 ciphertext (i.e. false negatives). All that the attack requires is that the attacker is capable of (correctly) identifying valid PKCS #1 v1.5 ciphertext sufficiently often.

D. Manger’s Attack

Following Bleichenbacher’s work, Manger [41] presented another padding oracle attack that allows an attacker to compute $c^d \bmod N$ without knowing the secret exponent d . While Manger’s attack was originally designed for attacking PKCS #1 v2.0 encryption, it can also be adapted to the PKCS #1 v1.5 case, obtaining better efficiency compared to Bleichenbacher’s attack, but requiring different prerequisites.

Attack Prerequisites. In this case we assume the existence of an oracle Ma which on input a ciphertext c answers whether $c^d \bmod N$ has the top byte equal to zero. More formally, let (N, d) be an RSA private key. The oracle Ma outputs the following for every ciphertext c

$$\text{Ma}(c) = \begin{cases} 1 & \text{if } c^d \bmod N \text{ starts with } 0x00 \\ 0 & \text{otherwise} \end{cases}.$$

That is, the oracle outputs for a given ciphertext c whether its decryption $c^d \bmod N$ lies in the interval $[0, B - 1]$ or not, where $B = 2^{8(\ell-1)}$ and ℓ is the byte-length of N .

High Level Attack Description. Let $c = m^e \bmod N$ be a ciphertext. At a high level, Manger’s attack is very similar to Bleichenbacher’s attack, requiring the attacker to choose a value s , to compute $c^* \leftarrow c \cdot s^e \bmod N$ and to query Ma in an attempt to find a c^* such that $\text{Ma}(c^*) = 1$.

Attack Efficiency. Compared to Bleichenbacher’s original attack which requires around 1 million oracle calls in order to perform an RSA secret operation, Manger’s attack only requires a little more than $\log_2(N)$ oracle calls. However, in contrast to Bleichenbacher’s attack which can tolerate oracle false negatives, Manger’s attack requires a “perfect” oracle which always answers correctly, without any errors.

E. The Interval Oracle Attack

Well before Bleichenbacher’s attack, Ben-Or et al. [8] proved the security of single RSA bits, by showing an algorithm for decrypting RSA ciphertexts given one-bit of plaintext

leakage. One of the oracles considered in that work is the interval oracle, that indicated if the plaintext is inside or outside a specific interval.

Attack Prerequisites. More specifically, for an RSA private key (N, d) assume that we have an oracle that outputs the following for every ciphertext c

$$\text{In}(c) = \begin{cases} 1 & \text{if } c^d \bmod N \text{ starts with bit } 1 \\ 0 & \text{otherwise} \end{cases}.$$

That is, the oracle outputs for a given ciphertext c whether its decryption $c^d \bmod N$ lies in the interval $[0, 2^{8\ell-1}]$ or not, where ℓ is the byte-length of N .

High Level Attack Description. The main idea behind the attack is to generate two random multiplications $c_1 = a \cdot c$ and $c_2 = b \cdot c$ of the ciphertext c , and then use an euclidean greatest common divisor (gcd) algorithm to compute $\text{gcd}(c_1, c_2)$. When a pair of ciphertext c_1, c_2 is found such that $\text{gcd}(c_1, c_2) = 1$, it is possible to efficiently recover $c^d \bmod N$. See [8] for a more complete description.

Attack Efficiency. The attack of [8] is relatively efficient, requiring about $15 \log_2 N$ oracle queries in order to decrypt the ciphertext c . For a random choice of c_1 and c_2 the attack will succeed with probability $6/\pi^2$.

F. Notation and Additional Padding Oracle Attacks

Following the attacks of Ben-Or et al. [8], Bleichenbacher [10], and Manger [41], several follow-up works obtained similar results using additional padding oracle commonly found in implementations of PKCS #1 v1.5, where some oracles provide more information than others [7, 36]. In this paper, we consider four different checks that an implementation can validate against the RSA-decrypted padded plaintext. All implementations start by checking that the padded plaintext starts with $0x0002$, and then may proceed with more checks.

- The first check corresponds to the test for a zero byte somewhere after the first ten bytes of the plaintext.
- The second check verifies that there are no zero bytes in the padding string PS'' .
- The third check verifies the plaintext length against some specific value (48 byte for a TLS premaster secret in our case).
- Finally, the fourth check is payload-aware and TLS-specific: it verifies the first two bytes of the payload; these bytes are set to the client’s protocol version as defined in [18].

Notation. Following the notation of Bardou et al. [7] we use FFFF to denote an oracle that gets as input a ciphertext and returns true in case the corresponding plaintext passes all four checks (and returns false otherwise). This is the same as the Bad-Version Oracle (BVO) from Klíma et al. [36]. Next, we use FFFT to denote an oracle that returns true in case the plaintext passes the first three checks, and false otherwise (disregarding the last check), FFTT to denote an oracle that returns true in case the plaintext passes the first two checks

(and disregards the last two checks), FTTT to denote an oracle that returns true in case the decrypted plaintext passes the first check (and disregards the latter three checks). Finally, we use TTTT to denote an oracle that disregards the four checks, returning true in case the plaintext corresponding to the input ciphertext starts with 0x0002.

Notice that in this notation, the Bleichenbacher oracle described in Section II-C corresponds to FTTT. Finally, we use M to denote Manger’s oracle presented in Section II-D and I to denote the Interval oracle presented in Section II-E.

G. The TLS Mitigation for the Bleichenbacher attack

The TLS specifications [16, 17, 18] define mitigations to implement in order to prevent the Bleichenbacher attack. More specifically: a server is always supposed to generate a random premaster secret; if the premaster secret decryption and PKCS #1 v1.5 validation fails, the server needs to use the randomly generated premaster secret to continue the TLS handshake. As this handshake will only fail later on, the attacker is deprived of the ability to tell if the modified ciphertext used during the attack was PKCS #1 v1.5 conforming or not. If the decryption is valid, the obtained premaster secret is used. This choice of premaster secret depending on the validity of the padding must be done in constant-time as well. Unfortunately, correctly implementing these mitigations is an extremely difficult task. Any differences in the server’s behavior between the PKCS #1 v1.5 conforming and the non-confirming cases may be exploited to obtain a Bleichenbacher-type oracle [11, 42]. Although most implementations have attempted to implement constant-time code for these mitigations, we show that all but two are still vulnerable to microarchitectural side-channel attacks.

H. Microarchitectural Side Channels

To improve the performance of programs, modern processors try to predict the future program behavior based on its past behavior. Thus, processors typically cache some microarchitectural state that depends on past behavior and subsequently use that state to optimize future behavior. Unfortunately, when multiple programs share the use of the same microarchitectural components, the behavior of one program may affect the performance of another, thereby creating contention on shared resources. Microarchitectural side channel attacks exploit this contention in order to leak otherwise unavailable information between programs [24].

Cache-Based Side Channel Attacks. Caching components, and in particular data and instruction caches, are often exploited for microarchitectural attacks. Cache-based attacks have been used to retrieve cryptographic keys [2, 9, 27, 32, 40, 48, 49, 56, 60], monitor keystrokes [28], perform website fingerprinting [47], and attack other algorithms [12, 58]. At a high level, cache attacks typically follow one of two patterns, which we now discuss.

FLUSH+RELOAD. In the FLUSH+RELOAD [60] attack and its variations [28, 29, 61], the attacker first evicts (flushes) a memory location from the cache. The attacker then waits a

bit, before reloading the flushed location, measuring the time to access the evicted memory. If the victim accesses the same memory location between the flush and the reload phases, the memory will be cached, and access will be fast. Otherwise, the memory location will not be cached and the access will be slower. Thus, the attacker deduces information regarding the victim’s access patterns to a given address.

PRIME+PROBE. Attacks employing the PRIME+PROBE technique [48, 49] or similar techniques [2, 19, 32, 40], first fill the cache with the attacker’s data; the attacker then yields its process, allowing the victim to execute code before measuring access-time to his own data. When the victim accesses its data, this data will evict some of the attacker’s data from the cache. By measuring the access time to the attacker’s data, the attacker can infer some information on the victim’s memory access patterns.

Attack Limitations. Both attacks require that the victim and attacker share some CPU caching components, implying that both programs have to run on the same physical machine. At a high level, FLUSH+RELOAD tends to be more accurate and have fewer false positives than PRIME+PROBE [60]. However, FLUSH+RELOAD requires the attacker to share memory with the victim and thus is typically applied to monitoring the victim’s code execution patterns, rather than data access patterns.

Branch-Prediction Attacks. The branch predictor of the processor has also been a target for microarchitectural attacks [3, 21, 22, 23, 38]. The branch predictor typically consists of two components, the Branch Target Buffer (BTB) which predicts branch destinations, and the Branch History Buffer (BHB), also known as the directional predictor, which predicts the outcome of conditional branches.

When a program executes a branch instruction, the processor observes the branch outcome and destination and modifies the state of the branch predictor. Attacks on the branch predictor exploit either the timing differences between correct or incorrect prediction or the performance monitoring information that the processor provides to recover the state of the predictor and detect the outcomes of prior branches executed by a victim program.

To mitigate Spectre attacks [37], Intel introduced mechanisms for controlling the branch predictor [31]. It is not clear whether these mechanisms completely eliminate branch prediction channels [25]. Furthermore, we have verified that by default Ubuntu Linux does not use the Indirect Branch Predictor Barrier mechanism to protect user processes from each other.

III. ATTACK MODEL AND METHODOLOGY

In this paper, we target implementations of PKCS #1 v1.5 that leak information via microarchitectural side channels. We then exploit the leaked information to implement a padding oracle, which we use to decrypt or to sign a message. To mount our attacks the adversary needs three capabilities:

1. Side Channel Capability. The first capability an adversary needs is to mount a microarchitectural side channel attack

against a vulnerable implementation. For that, the adversary needs the ability to execute code on the machine that runs the victim’s implementation.

2. Privileged Network Position Capability. Our attacks exploit a padding oracle attack to perform a private key operation such as a signature or decryption of a message that has been sent to the victim. In order to decrypt a ciphertext and use its result, an adversary must first obtain a network man-in-the-middle position. In order to forge signatures, an adversary must first obtain the relevant data to sign and be in a privileged position to exploit it.

3. Decryption Capability. The third capability our adversary needs is the ability to trigger the victim server to decrypt ciphertexts chosen by the adversary.

A concrete attack scenario we consider in this work is attacking a TLS server running on the same physical hardware as an unprivileged attacker. For example, a TLS server running in a virtual machine on a public cloud server, where the physical server hardware is shared between the victim’s TLS server and an attacker’s virtual machine. Indeed, previous works have shown that attackers can achieve co-location [50], and leverage it for mounting side channel attacks [30]. Thus, the first capability is achievable for a determined adversary.

The second and third capabilities are achievable in this scenario by an attacker that controls any node along the path between the client and the server. Malicious network operators are one example of actors that have such control, but this is not the only case. In particular, attackers can exploit vulnerabilities in routers to assume control and mount our attack[15].

This scenario, however, also presents specific problems; TLS 1.3 does not support RSA key exchanges anymore and in TLS 1.2 (Elliptic Curve) Diffie-Hellman key exchanges are recommended over RSA key exchanges. Hence, the adversary needs to perform active protocol downgrade attacks to force the use of RSA in the communication. Furthermore, the adversary may need to perform a large number of decryptions, which is time consuming and is limited to each client’s (e.g. browsers) timeout configurations. We explain in [Section VI](#) how it is possible to perform such downgrade attacks.

IV. VULNERABILITY CLASSIFICATION

We now present the code of several popular implementations of RSA PKCS #1 v1.5 decryption and discuss their vulnerabilities to microarchitectural padding oracle attacks. Before describing the different types of vulnerabilities we discovered, we provide an outline of typical TLS implementations.

- **Data Conversion.** First, the RSA ciphertext is decrypted and the resulting plaintext is converted into a byte array.
- **PKCS #1 v1.5 Verification.** Next, the conformity of the array to the PKCS #1 v1.5 standard is checked.
- **Padding Oracle Mitigations.** Finally, in case the array is not PKCS #1 v1.5 conforming, the server must somehow abort the protocol without revealing this fact to the attacker (e.g., in the form of an error message). To achieve this,

```
1 int RSA_padding_check_none(to, tlen, from, flen){
2 // to is the output buffer of maximum length tlen
  bytes
3 // from is the input buffer of length flen bytes
4 memset(to, 0, tlen - flen);
5 memcpy(to + tlen - flen, from, flen);
6 return tlen;
7 }
```

Listing 1. Pseudocode of raw plaintext copy with no padding check function

the server executes the TLS padding oracle attacks countermeasure, and proceeds with the TLS handshake using random data. We note here that any code executing prior to completing this stage is very sensitive, as any information leakage regarding the PKCS #1 v1.5 compliance of the byte array might be leveraged to form a padding oracle.

Unfortunately, despite more than twenty years of research in both padding oracle attacks and side channel resistance, in this work we find that vulnerabilities still occur in all of these stages. We now provide a high level description of the various stages and their associated side channel vulnerabilities.

A. Data Conversion.

In RSA, the plaintext and the ciphertext are large numbers, e.g. 2048-bit long. These are typically represented as little-endian arrays of 32- or 64-bit words. PKCS #1 v1.5, however uses big-endian byte arrays, thus requiring a format conversation. For values of fixed length, this conversation is relatively straightforward. However, while the length of the RSA modulus provides an upper bound on the length of the RSA decryption result, the exact length of the RSA plaintext is not known until after RSA decryption of the corresponding ciphertext. Thus, in case the RSA decryption result is too short, the little-to-big endian conversation code has to pad the ciphertext with a sufficient amount of zero bytes.

As an example, consider the pseudo code of the implementation of the OpenSSL function `RSA_padding_check_none` in [Listing 1](#). The function is called as part of the implementation of the TLS protocol in OpenSSL, and its purpose is to copy the RSA decryption results to an output buffer, without performing any padding checks.

To handle the case that the plaintext from the RSA decryption is smaller than the output buffer, `RSA_padding_check_none` uses `memset` to pad the output buffer. The length of the padding is set to the difference between the lengths of the output array and the plaintext. In case of a full-length plaintext, the length of the padding is zero. Using Branch-Prediction attack we can detect this scenario, and learn whether the plaintext is full-length. This gives us the oracle required for a Manger attack.

Unfortunately, this example is by no means unique, and multiple implementation expose FTTT- or Manger-type padding oracles during the data conversion phase. See [Appendix A](#) for further examples.

B. PKCS #1 v1.5 Verification

Once the data is represented as a sequence of bytes, the implementation needs to check that it is PKCS #1 v1.5 con-

forming, that is, that the first byte is zero, the second is 0x02, the following eight bytes are non-zero, and that there is a zero byte at a position above 10. Yet, many implementations branch on the results of these checks, leaking the outcome to a side channel attacker via the implementation’s control flow. The exact oracle obtained depends on the specific implementation and the type of leakage.

The OpenSSL RSA PKCS #1 v1.5 decryption API provides an example of such an issue. OpenSSL is the cryptographic engine used in many applications. It exports a function, `RSA_public_decrypt`, whose arguments are an input buffer containing the ciphertext, an output buffer for the plaintext, the RSA decryption key, and the padding mode to check the plaintext against. When using PKCS #1 v1.5 padding, `RSA_public_decrypt` invokes `RSA_padding_check_PKCS1_type_2` to validate the padding after decryption. A pseudocode of the function is shown in Listing 2.

As the pseudocode shows, OpenSSL performs the checks outlined in Section II-B in constant-time (Lines 7–13), returning the length of the decrypted message if the decryption is successful, and `-1` if there is a padding error. To set the return value, the function uses an explicit branch (Line 17). Furthermore, the memory copy in Line 21 is only executed in case of a successful decryption, whereas the error logging (Line 25) is invoked in the case of a padding error.

A comment in the code (Line 15) indicates that the authors are aware of the leakage, and the manual page for the function warns against its use [46]. Thus, OpenSSL does not use this PKCS #1 v1.5 verification code for its own implementation of the TLS protocol. Furthermore, both Xiao et al. [57] and Zhang et al. [62] exploit the leakage through the conditional error logging for mounting Bleichenbacher attacks.

Amazon’s s2n. OpenSSL is the cryptographic engine underlying many applications, all of these are potentially vulnerable to our cache-based padding oracle attack. Specifically, Amazon’s implementation of the TLS protocol, s2n [54], uses this API, and consequently leaks an FFFT-type oracle. For other vulnerabilities in s2n, see Appendix A-B.

C. Padding Oracle Mitigations.

As Section II-G describes, when a TLS implementation detects that a plaintext does not conform to the PKCS #1 v1.5 format, it cannot just terminate the handshake, because this creates a padding oracle. Instead it must replace the non-conforming plaintext with a random sequence of bytes and proceed with the TLS handshake. However, some implementation fail to protect this replacement, leaking the deployment of the countermeasure and allowing the creation of a padding oracle.

We can find examples of such leakage in CoreTLS, Apple’s implementation of the TLS protocol that is sometimes used in both MacOS and iOS devices.

Listing 3 shows the code that handles the mitigation of Bleichenbacher’s attack in CoreTLS (i.e., replacing the incorrectly-padded RSA plaintext with random data). Lines 7 and 8 perform the RSA decryption and the validation of the PKCS #1 v1.5 format. The code logs validation failure of

TABLE I
SUMMARY OF IDENTIFIED PADDING ORACLES.

	Data Conv.	PKCS #1 v1.5 Verification	TLS Mitigation
OpenSSL	M	M	
OpenSSL API	M	FFTT	
Amazon s2n		FFFT	
MbedTLS	I	FFTT, FFFT*	
Apple CoreTLS			FFTT, FFFT, FFFF
Mozilla NSS	M	M, TTTT, FTTT*	FFFF
WolfSSL	M	M, FFFT	FFTT, FFFF
GnuTLS	M	M, TTTT, FTTT	FFTT, FFFT
BoringSSL		<i>Not Vulnerable</i>	
BearSSL		<i>Not Vulnerable</i>	

the PKCS #1 v1.5 format in Line 11. It also checks that the output is of the expected length, issuing a log message on failure (Lines 13–17). For brevity we omit the code that handles the success case (Line 20). The main mitigation against Bleichenbacher attacks occurs in Line 24, where the code generates a random value to be used as the session key.

While CoreTLS’ code that verifies the PKCS #1 v1.5 padding is constant time, the code that handles the mitigations against padding oracle attacks (Listing 3) is far from constant time. As seen in Listing 3, the code contains multiple sources of side channel leakage which we now describe.

First, all of the conditional if statements in the presented code can be exploited by branch prediction attacks to implement FFFT (Line 9), FFFT (Line 13), or FFFF (Lines 19 and 22) Bleichenbacher-type oracles.

Next, a cache attack can monitor either the code of the log message function or the code of the random number generator, which only runs if the PKCS #1 v1.5 validation fails. Another option is to monitor the bodies of the if statements in Lines 19 or 22. These attacks be used to implement an FFFF-type padding oracle.

Finally, generating the random session key only on PKCS #1 v1.5 validation failure (Line 24) is a significant weakness in the implementation. Random number generation is a non-trivial operation that may take significant time and thus might expose a Bleichenbacher oracle via a timing attack. That is, by simply measuring the response time of a TLS server that uses the CoreTLS library, an attacker might get a FFFF-type Bleichenbacher oracle.

D. Summary of the Findings.

Table I summarizes our findings, showing the identified oracle types in each stage of the implementations we evaluated. As we can see, seven of the nine tested implementations expose padding oracles via microarchitectural attacks. Only BearSSL and Google’s BoringSSL are not vulnerable to our attacks.

V. EXPERIMENTAL RESULTS

To validate that the vulnerabilities we identified can indeed be exploited, we mounted concrete side-channel attacks on some of the implementations. We now discuss some of the techniques we used for this validation.

```

1 int RSA_padding_check_PKCS1_type_2(to, tlen, from, flen, num_bytes){
2 // to is the output buffer of maximum length tlen bytes
3 // from is the input buffer of length flen bytes
4 // num_bytes is the maximum number of bytes in an RSA plaintext
5 // returns the number of message bytes (not counting the padding) or -1 in case of a padding error
6
7 good = constant_time_is_zero(from[0]);
8 good &= constant_time_eq(from[1], 2);
9 zero_index = find_index_of_first_zero_byte_constant_time(from+2, flen);
10 good &= constant_time_greaterOrEqual(zero_index, 2 + 8); //first 10 plaintext bytes must be non-zero
11 msg_index = zero_index + 1; //compute location of first message byte
12 msg_len = num_bytes - msg_index; //compute message length
13 good &= constant_time_greaterOrEqual(tlen, msg_len); //check that to buffer is long enough
14
15 /* We can't continue in constant-time because we need to copy the result and we cannot fake its length.
16    This unavoidably leaks timing information at the API boundary. */
17
18 if (!good) {
19     mlen = -1;
20     goto err;
21 }
22 memcpy(to, from+msg_index, mlen);
23
24 err:
25 if (mlen == -1)
26     RSAerr(RSA_F_RSA_PADDING_CHECK_PKCS1_TYPE_2, RSA_R_PKCS_DECODING_ERROR);
27 return mlen;
28 }

```

Listing 2. Pseudocode of RSA_padding_check_PKCS1_type_2

```

1 int SSLDecodeRSAKeyExchange(keyExchange, ctx){
2 keyRef = ctx->signingPrivKeyRef;
3 src = keyExchange.data;
4 localKeyModulusLen = keyExchange.length;
5 ... // additional initialization code omitted
6
7 err = sslRsaDecrypt(keyRef, src,
8 localKeyModulusLen,
9 ctx->preMasterSecret.data,
10 SSL_RSA_PREMASTER_SECRET_SIZE, &outputLen);
11 if(err != errSSLSuccess) {
12     /* possible Bleichenbacher attack */
13     sslLogNegotiateDebug("SSLDecodeRSAKeyExchange:
14 RSA decrypt fail");
15 } else if(outputLen !=
16 SSL_RSA_PREMASTER_SECRET_SIZE) {
17     sslLogNegotiateDebug("SSLDecodeRSAKeyExchange:
18 premaster secret size error");
19     // not passed back to caller
20     err = errSSLProtocol;
21 }
22 if(err == errSSLSuccess) {
23     ... // (omitted for brevity)
24 }
25 if(err != errSSLSuccess) {
26     ... // (omitted for brevity)
27     sslRand(&tmpBuf);
28 }
29 /* in any case, save premaster secret (good or
30 bogus) and proceed */
31 return errSSLSuccess;
32 }

```

Listing 3. Apple’s TLS mitigation function

A. Attacking the OpenSSL API

The vulnerability in the OpenSSL API (Section IV-B) has already been disclosed by both Xiao et al. [57] and Zhang et al. [62]. Our attack is similar to the attack of Zhang et al. [62], but achieves a significantly lower error rate, resulting

in a lower number of required oracle invocations. Combined with our improved error handling (Section VI-B) we achieve a reduction by a factor of 6 in the number of oracle queries we require.

Our test machine uses a 4 core Intel Core i7-7500 processor, with a 4 MiB cache and 16 GiB memory, running Ubuntu 18.04.1. We use the Flush+Reload attack [60], as implemented in the Mastik [59] toolkit.

To reduce the likelihood of errors, we monitor both the call-site to RSAerr (Line 25 of Listing 2) and the code of the function RSAerr. Monitoring each of these locations may generate false positives, i.e. indicate access when the plaintext is PKCS #1 v1.5 conforming. The former results in false positives because the call to RSAerr shares the cache line with the surrounding code, that is always invoked. The latter results in false positives when unrelated code logs an error. By only predicting a non-conforming plaintext if both locations within a short interval, we reduce the likelihood of false positives.

We note that this technique is very different to the approach due to Genkin et al. [26] of monitoring two memory locations to reduce false negative errors due to a race between the victim and the attacker [6]. Unlike us, they assume access if any of the monitored locations is accessed.

Overall, our technique achieves a false positive rate of 4.3% and false negative rate of 1.1%.

B. Attacking the OpenSSL Data Conversion

We now turn our attention to the code OpenSSL uses for its own implementation of the TLS protocol. As discussed in Section IV-A, OpenSSL leaks a Manger oracle through the length argument in the call to memset in Line 4 of Listing 1.

We now show how we detect that the length passed to `memset` is zero.

We implement a proof-of-concept attack on an Intel NUC computer, featuring an Intel Core i7-6770HQ CPU, with 32 GiB memory, running Centos 7.4.1708. The GNU C library provides multiple implementations for `memset`, each optimized for different processor features. During initialization, the library chooses the best implementation for the computer, and stores it in a function pointer. In run time, the program invokes the best implementation of `memset` by dereferencing the function pointer. On our system, the selected function is `__memset_sse2`. We show part of the (disassembled) code of this function in Listing 4.

```
1 <+209>: test    $0x1,%dl
2 <+212>: je     0x40e918 <__memset_sse2+216>
3 <+214>: mov    %cl, (%rdi)
4 <+216>: test    $0x2,%dl
5 <+219>: je     0x40e87a <__memset_sse2+58>
6 <+225>: mov    %cx, -0x2(%rax,%rdx,1)
7 <+230>: retq
```

Listing 4. A snippet of `__memset_sse2`

The presented code is only executed if the length argument for `memset` is less than 4. Line 1 of the code first test the least significant bit of the length, branching, in Line 2 if it is clear, i.e. if the length is 0 or 2. In Line 4, the code tests if the second bit of the length, branching in Line 5 if the length is less than 2. Thus, if both branches at Lines 2 and 4 are taken, the length argument is 0.

Branch Prediction Attack. Our attack follows previous works in creating *shadow* branches, at addresses that match the least significant bits of monitored branches [22, 38]. Because the branch predictor ignores the high bits of the address, the outcome of the victim branch affects the prediction for the matching shadow branch. That is, when a monitored branch is taken, the BTB predicts that both the monitored branch and its shadow will branch to the same offset as the monitored branch.

Prior works either measure the time to execute the shadow branch [22] or check the performance counters [38] to detect mispredictions of the shadow branch, and from these infer the outcome of previous executions of the monitored branch. However, performance counters are not always available to user processes, and measurements of execution time of branches are noisy. Instead, we combine the branch prediction attack with FLUSH+RELOAD [60] to achieve high accuracy detection of mispredictions.

Specifically, for each monitored branch we create two shadows, the *trainer* and the *spy* branches. Each of these branches to a different offset, such that the offsets of the monitored branch and of the shadow branches each falls in a different cache line. The attack then follows a sequence of steps:

- Invoke the trainer shadow to train the branch predictor to predict the trainer offset for all three branches.
- Flush the cache line at the trainer offset from the spy branch from the cache.

- Execute the victim. If the victim branch is taken, it will update the BTB state to predict the victim offset for all three branches
- Invoke the spy branch. Because the branch predictor predicts either the victim or the trainer offset, the spy branch mispredicts. In the case that the victim branch has not been taken, the mispredicted branch will attempt to branch to the trainer offset from the branch, bringing the previously flushed line back into the cache.
- Measure the time to access the previously flushed line. If the victim branch has been taken, this line will not be cached, and access will be slow. If, however, the victim branch did not execute or was not taken, the line will be in the cache due to the misprediction in the previous step, and access will be fast.

We implemented this attack and we can predict the outcome of each of the monitored branches with a probability higher than 98%. We cannot, however, monitor both branches concurrently. Consequently, for the manger attack, we will have to send each message twice. Once for monitoring the outcome of the branch in Line 2 and the other for the branch in Line 5.

VI. MAN IN THE MIDDLE ATTACKS

In principle, a padding oracle attack can be used to compute the premaster secret used in any TLS connection that uses RSA key exchange, thereby completely breaking the connection’s confidentiality guarantees. At the time of writing, this accounts for $\approx 6\%$ of all TLS connections[1, 45]. Moreover, the attack does not need to be performed online, as an attacker can record the connection’s encrypted traffic and use a padding oracle to decrypt it at a later date. Even though the attack seems to have only limited applicability, we show how to use an online downgrade attack to break the security of all TLS connections, even when they do not wish to use the RSA option.

Man In The Middle Attacks via Padding Oracles. Next, as noted in [35], efficient padding oracle attacks can be used *online* to mount man in the middle (MiTM) downgrade attacks on TLS connections. Such online padding oracle attacks are particularly dangerous as an attacker mounting a downgrade attack can force both the client and server to initiate the handshake with TLS 1.2 RSA key exchange, break the connection’s premaster secret using the padding oracle, and subsequently finish the handshake using the obtained premaster secret.

Moreover, efficient online padding oracle attacks are dangerous even in the case where an updated client uses a protocol that does not support RSA key exchange (such as TLS 1.3 and Quick) to connect to a server that supports RSA key exchange. Assuming the server uses the same certificate for RSA decryption and signing, an attacker can simulate a TLS 1.3 connection with the client, and use the online padding oracle to sign a forged ephemeral public key (see [35] for details). As observed by Jager et al. [35] using single-certificate scenario for all RSA operations across all protocols in a given server is in fact quite common, as popular TLS

servers often do not use separate certificates for RSA key exchange and RSA signing¹.

As TLS servers often offer older protocols which contain padding oracles for compatibility reasons, and will continue to do so for many years, the above-outlined attack scenarios makes padding oracles a threat to *almost all* TLS connections.

Mounting Online Padding Oracle Attacks. Mounting such an online padding oracle MiTM attack, the attacker has to finish the attack before the browser-enforced TLS timeout. As was shown by [4], we can cause Firefox to keep a TLS handshake alive indefinitely, thus allowing us to perform even very long attacks. Using a BEAST style attack [20] we can perform this attack in the background, without the user noticing any long delays. However, to mount such an online padding oracle MiTM attack against other browsers, the attacker has to be extremely efficient, finishing the attack before the browser-enforced TLS timeout. Unfortunately, this is often difficult to achieve as a padding oracle might require the attacker to perform many thousands of TLS handshakes with the server, which takes much longer than a typical browser-enforced TLS timeout of about 30 seconds (for Google Chrome and Microsoft Edge).

Analysis and improvement of Padding Oracle Attacks. In this section, we analyze the complexity of padding oracle attacks for an online MiTM scenario. Our contributions are as follows. First, we present a novel analysis of the query complexity required from a padding oracle attack (Section VI-A). Next, we handle the case of imperfect and noisy oracles (Section VI-B). Finally, in Section VII we address the question of parallelizing padding oracle attacks across any available number of servers, demonstrating a new connection between padding oracle attacks and lattice reduction techniques.

A. Reducing the Query Complexity of Padding Oracles

Assume we would like to break the security of a specific account in some popular online service (e.g., Gmail). As the connection is usually done via https (which uses TLS), one attack vector is to attempt to break the user’s existing TLS connection with the online service. Using padding oracles to mount a MiTM downgrade attack on a specific connection might be difficult given the 30 seconds browser-enforced timeout for completing the TLS handshake. In our new analysis, we assume that we perform a BEAST style attack [20]. In this scenario a malicious web site controlled by the attacker, causes the user’s browser to repeatedly try to connect to the TLS server in the background without the user’s knowledge. This attack only requires that the browser supports JavaScript, and does not need any special privileges (in particular, the attacker does not have to compromise the normal operation of the target machine in any way). A successful MiTM attack on even a single TLS handshake will allow the attacker to decrypt the user’s login token, thereby allowing a malicious server login.

¹For example, at the time of writing, Amazon AWS servers uses the same RSA certificate for signing and key encryption.

Low Success Probability is Sufficient. Any Bleichenbacher like attack requires a large expected number of queries that can not be completed before the browser’s timeout. However, we can use the long tail distribution of the number of queries in order to find an attack that required a much lower number of queries with some small but not negligible probability (say 1/1000). By using the BEAST attack we can amplify this small probability by forcing the browser to repeatedly negotiating new TLS handshakes in the background until we succeed.

The remainder of this section is as follows. First, following the outline of Bleichenbacher’s and Manger’s attacks presented in Sections II-C and II-D, we notice that the overall oracle query complexity of both attacks highly depends on the probability that the padding oracle outputs 1 on a random ciphertext c . We thus begin by analyzing this probability for several common oracles. Next, we present simulations results of the total oracle query complexity required from several attack success probabilities.

Analyzing OpenSSL API FFFT Oracle. We begin by analyzing the probability that a random ciphertext c conforms with the requirements of the FFFT padding oracle present in OpenSSL’s decryption API (Section IV-B). Let (d, N) be an RSA private key. For a ciphertext c to be conforming, the following must hold:

- 1) First, the two topmost bytes of $c^d \bmod N$ (the RSA plaintext corresponding to c) must be 0x0002. For a random c , this happens with probability of 2^{-16} .
- 2) Next, the first 8 padding bytes of the plaintext corresponding to c must be non-zero. For a random c , this event happens with probability of $(255/256)^8$.
- 3) The plaintext corresponding to c contains a zero byte. For a 2048-bit RSA modulus N , we have 246 remaining bytes. Thus, for a randomly selected c , this event holds with probability of $1 - (255/256)^{246}$ (or $1 - (255/256)^{502}$ for a 4096-bit modulus).

We obtain that for any 2048-bit RSA private key it holds that

$$\begin{aligned} \Pr_c[\text{FFTT}(c) = 1] &= 2^{-16} \cdot \left(\frac{255}{256}\right)^8 \cdot \left(1 - \left(\frac{255}{256}\right)^{246}\right) \\ &\approx 9.14 \cdot 10^{-6}. \end{aligned}$$

Similarly, for any 4096-bit RSA private key, we obtain that $\Pr_c[\text{FFTT}(c) = 1] \approx 1.27 \cdot 10^{-5}$. Next, the expected number of oracle queries required to obtain a conforming ciphertext is $1/\Pr_c[\text{FFTT}(c) = 1]$ which results in about 110000 queries for 2048-bit key and about 80000 queries for 4096-bit key. However, the median (50% success probability to obtain a conforming ciphertext) is about 76000 queries for 2048-bit key and about 55000 queries for 4096-bit key. Next, if we allow a low success probability of 1/1000, the required number of random oracle queries to obtain a conforming ciphertext is much lower, approximately 110 queries for 2048-bit key and about 80 queries for 4096-bit key. Finally, we note that partially due to this effect, the total query complexity of our attack goes down as the RSA private key gets larger.

Analyzing MbedTLS FFFT Oracle. We now proceed to analyze the FFFT padding oracle present in MbedTLS implementation of the PKCS #1 v1.5 verification code (Section A-C). Let (d, N) be an RSA private key. For a plaintext c to be conforming to an FFFT oracle, the following must hold.

- 1) The first two conditions of the FFFT oracle present in the OpenSSL decryption API hold. For a random ciphertext the probability that both conditions hold is $2^{-16} \cdot (255/256)^8$.
- 2) The size of the unpadded plaintext corresponding to c is between 0 and 48 bytes. For a 2048-bit RSA key, we have 256 bytes of padded plaintext. The first 10 bytes are checked in the first condition, leaving 246 bytes for the padding and the plaintext itself. As the padding string must consist of some number of non zero bytes and terminate with a zero byte, we obtain that for a random 2048-bit ciphertext c , this event holds with probability of $(255/256)^{246-48} \cdot (1 - (255/256)^{48})$. Similarly, for 4096-bit RSA key (containing 512 bytes), this event holds for a random ciphertext with probability of $(255/256)^{502-48} \cdot (1 - (255/256)^{48})$.

Thus, for any 2048-bit RSA private key it holds that

$$\Pr_c[\text{FFFT}(c) = 1] = \frac{1}{2^{16}} \left(\frac{255}{256}\right)^8 \left(\frac{255}{256}\right)^{198} \left(1 - \left(\frac{255}{256}\right)^{48}\right) \approx 1.16 \cdot 10^{-6}.$$

Similarly, for any 4096-bit RSA private key, we obtain that $\Pr_c[\text{FFFT}(c) = 1] \approx 4.28 \cdot 10^{-7}$. Next, the expected number of oracle queries required to obtain a conforming ciphertext is $1/\Pr_c[\text{FFFT}(c) = 1]$ which results in about 857000 queries for 2048-bit key and about 2334000 queries for 4096-bit key. As in the case of the FFFT analyzed previously the query complexity is greatly reduced for smaller success probabilities, requiring about 594000 queries for a 50% success probability and only 860 queries for a success probability of 1/1000 for 2048-bit keys. For 4096-bit keys, the median is 1618000 queries and 2300 queries are required for a success probability of 1/1000.

Analyzing the Manger Oracle. The Manger attack complexity is much simpler, having the number of queries required be approximately the length of the RSA modulus in bits with very low variance (i.e., about 2048 queries for 2048-bit keys and 4096 queries for 4096-bit keys).

Full Attack Simulation. While the query complexity of the entire padding oracle attack highly depends on the probability p that the padding oracle outputs 1 on a random ciphertext, for Bleichenbacher-type oracles the exact relation between p and the attacks' query complexity is rather difficult to analyze. Instead, we ran 500000 simulations of the full attack using the FFFT, FFFT and Manger type oracles, for a 2048-bit RSA modulus. The results of our simulation are presented in Table II, for both decryption and signature forging attacks. For each oracle type and attack type, we give the required number of oracle queries needed to complete the attack with the different success probabilities. It is easy to see that for

a small success probability (i.e., 1/1000), we need a much lower (by up to a factor of 10) number of padding oracle queries, compared to the number required for a 50% success probability. As outlined above, while the success probability of each individual attack attempt is low (1/1000), the attacker can always use BEAST-style techniques, having a malicious website repeatedly issue TLS connections to the target website. As soon as a single connection attempt is broken, the attacker can decrypt the user's login token, compromising the account. Finally, we note the because each attack attempt has a low oracle query complexity, it is possible to complete the attempt below the 30 seconds timeout enforced by Chrome and Edge.

B. Handling Oracle Errors

In case an implementation reveals a padding oracle directly via network messages, the resulting padding oracle is often "perfect", exhibiting no errors. However, in case the oracle is obtained via side channels (such as the microarchitectural padding oracles considered in this work) the result is often noisy, containing both false negative or false positive errors. Moreover, the error probability is often not symmetric, that is $\Pr[\text{False Positive}] \neq \Pr[\text{False Negative}]$. In this section we analyze the padding oracles considered in this work, presenting efficient strategies for error recovery. See Table II for a summary of the results.

Handling Errors in Manger Type Attack. As outlined in Section II-D the Manger attack is not error tolerant, having any type of error in any oracle query result in the attack failing to break the target TLS connection. Thus, to obtain an error-free result we propose to repeat each oracle query several times, taking a majority vote in the result. We now proceed to analyze the exact number of repetitions required by this approach.

Indeed, assume we want a padding oracle attack to succeed with a low probability of $p = 0.001$. For a 2048-bit RSA modulus, we will require about 2048 queries to break the target connection. This means that we require $(1 - \Pr[\text{error}])^{2048} > 0.001$ which yields $\Pr[\text{error}] < 1 - \sqrt[2048]{0.001} \approx 0.00337$. Next, from the experimental results outlined in Section V-B, we have that our side channel based Manger oracle has a false positive rate of 0.02 and a false negative rate of 0.02, as the Manger oracle cannot tolerate both types of errors, we obtain that each individual oracle calls is incorrect with a probability of at most 0.02. Assume we take the majority over r distinct oracle calls. For the majority to be incorrect, it has to be the case that $\Pr[\text{error}] < \sum_{i=r/2+1}^r (0.02)^i \approx (0.02)^{r/2+1} < 0.00337$, which yields $r = 3$.

Handling Errors in Bleichenbacher-type Oracles. As outlined in Section II-C, Bleichenbacher-type oracles are tolerant to one sided errors. Although a single false positive might cause the attack to fail, the attack can tolerate an arbitrary number of false negatives. This is since inserting a false negative (meaning the oracle answered 0 for ciphertext whose plaintext is PKCS #1 v1.5 conforming) into the attack algorithm will just require more queries to find the next blinded ciphertext that decrypts to a PKCS #1 v1.5 conforming plaintext. Although a false negative in the first phase can result

TABLE II
NUMBER OF ORACLE QUERIES REQUIRED FOR 2048-BIT RSA MODULUS.

Oracle	Signature Forging with Success Probability				Decryption with Success Probability			
	0.001	0.01	0.1	0.5	0.001	0.01	0.1	0.5
FFFT Oracle (OpenSSL API)	16381	19899	40945	122377	14700	15147	16764	50766
FFFT Oracle (MbedTLS)	139426	192633	533840	1292250	116699	123359	237702	870664
Manger Oracle	≈ 2048	≈ 2048	≈ 2048	≈ 2048	≈ 2048	≈ 2048	≈ 2048	≈ 2048
FFFT Oracle With Errors	29989	33944	57130	147406	28170	28683	30494	70990
Manger Oracle With Errors	≈ 6144	≈ 6144	≈ 6144	≈ 6144	≈ 6144	≈ 6144	≈ 6144	≈ 6144

in a large number of extra queries (the probability of finding a conforming plaintext is low), a false positive in the later stages might be relatively “cheap”.

To better understand the total query complexity required for a side-channel based Bleichenbacher-type oracle, we simulated the end-to-end attack using the false negative and false positive rates obtained in Section V-A (i.e., we set $\Pr[\text{False Positive}] = 0.043$ and $\Pr[\text{False Negative}] = 0.011$). As a Bleichenbacher-type oracle can tolerate false negatives, we used just one sided error correction. That is, in case the oracle reports the ciphertext as not PKCS #1 v1.5 compatible we accept the answer and try another ciphertext. However, in case the oracle reports that ciphertext does conform to the PKCS #1 v1.5 standard, we repeat the measurement 6 times (by issuing additional queries) and require that positive answer will be given for 5 out of the 6 queries. We note that this amount of repetitions was empirically chosen to minimize the attacks’ total query complexity. The results of our simulations can be seen in Table II. Notice that the total query complexity is between about $\times 1.2$ and about $\times 2$ the oracle queries required from the perfect oracle case (i.e., no false negatives or false positives).

VII. PARALLELIZATION OF THE ATTACK

We would like to exploit the fact that large service providers often use multiple TLS servers that share the same RSA key; this is done in order to support greater bandwidth, provide redundancy, and provide better latency for different locations around the world. We would like to be able to parallelize our attacks using several such servers as oracles, making the attack fast enough for a MiTM downgrade attack. Parallelization of the Bleichenbacher attack was first mentioned by Klíma et al. [36], who proposed to parallelize phase 2 of the attack Section II-C in case of multiple possible ranges, in a round robin manner on a single available server. Böck et al. [11] mentioned the possibility of using multiple servers to parallelize the attack. However, they have not given a concrete method of doing this, and did not address the inherent limitation of trivial parallelization methods.

Limitations of Trivial Parallelization. A trivial parallelization method for the Bleichenbacher attack is to parallelize multiple queries with different values for s_i in each phase of the attack; this will allow us to find the correct value faster. We can also use the approach of Klíma et al. [36] if we have multiple ranges in phase 2. Another approach is to parallelize the multiple identical queries for error cor-

rection in Bleichenbacher and Manger attacks mentioned in Section VI-B. However, both Bleichenbacher and Manger attacks are adaptive chosen ciphertext attacks, requiring at least $\log_2 N$ sequential queries. This is true even if we have an unlimited number of oracles.

Parallelization via the Closest Vector Problem. A new type of parallelization is possible by representing the problem as a variant of the Closest Vector Problem (CVP) and solving it with the LLL efficient lattice reduction technique [39]. Our main observation is that our adaptive attack can be parallelized by creating multiple parallel attacks against copies of the same ciphertext starting from different initial blinding values. We can then stop each attack after a predetermined number of adaptive queries, and combine the partial information recovered in each attack into a full plaintext recovery by using CVP.

A. Parallelization of the Manger Attack.

At each point in the Manger attack (after the initial blinding in phase 1), we know that $m \cdot s \bmod N$ is inside the interval $[a, b]$, where m is the unknown plaintext, s is the known blinding value, N is the RSA modulus, and the attack’s goal is to decrease the size of the interval $[a, b]$. During each adaptive step of the attack, the size of the interval is reduced. When $|[a, b]| = 1$, we know that $a = m \cdot s \bmod N$, and can recover the original plaintext by calculating $m = a \cdot s^{-1} \bmod N$. If we can approximate halve the size of the interval in each step, we can complete the attack after $\approx \log_2 N$ adaptive queries.

We look at the scenario where we run k attacks in parallel, but only have time for i adaptive queries for each attack. For each attack j after i queries the interval is $[a_j^i, b_j^i]$. The number of information bits we learned about the value of $m \cdot s \bmod N$ from the attack so far is given by $I_j^i = \log_2 N - \log_2 (b_j^i - a_j^i)$. If after i adaptive queries $\sum_{j=1}^k I_j^i > \log_2 N$ we can recover the value of m . We write k linear equations of the form

$$m \cdot s_j - a_j^i \bmod N < 2^{\log_2 (b_j^i - a_j^i)}$$

and recover m by solving CVP using an LLL lattice reduction algorithm. Full details about how to set up the lattice and how to interpret the short vectors that the LLL algorithm produces can be found in the full version of the paper.

Analyzing the Parallel Attack. We would like to analyze the trade-off between the number of adaptive queries and the number of parallel oracles. In the Manger attack the blinding phase requires on average 128 parallel queries, and gives us 8 bits of information on the plaintext. The next two phases (called step 1 and 2 in the original paper) are harder

to analyze, but experiments show that they usually require 40 – 100 adaptive queries and give us 8 – 12 extra bits of information. After that, each adaptive query gives us ≈ 1 bit of information. For RSA modulus of 2048 bits the original Manger attack without blinding requires 2100 adaptive queries and just one oracle (which requires negligible computation). On the other extreme we can try a fully parallelized attack using only the blinding phase. This will require approximately $128 \cdot 256 = 32768$ parallel queries, that will result in 256 equations giving us 8 bits each. Recovering the plaintext will require us to reduce a relatively large lattice of size ≈ 256 , which requires a considerable amount of computation. A more efficient trade-off will be to run several partial adaptive attacks in parallel.

Parallel Manger Attack Simulation. We ran a simulation to test the feasibility of performing a MiTM on a TLS connection and a 2048 bits RSA with multiple parallel partial Manger attacks. We assume that we have 30 seconds before the TLS connection will timeout and that each TLS handshake takes about 0.05 seconds (which is the actual time measured on a Core i7-7500U CPU @ 2.70GHz). We allow each of the parallel attacks to have 560 adaptive oracle queries (leaving us 2 seconds for the lattice reduction and finalizing the handshake). We simulated a parallel attack using 5 servers (the minimal number of servers required to fit at least 2048 queries in 30 seconds is 4, but due to overheads we require at least 5 servers).

We start by running the blinding phase in parallel until we get 5 valid blinding values. We then use our remaining queries to continue the 5 attacks in parallel. As before, we ran a simulation of the attack with 5 oracles 500000 times. With probability 0.001 we got at least 438 bits of information from each of the 5 attacks, or a total of more than 2190 bits. This is more than the required number of bits to recover the plaintext. We successfully implemented and tested a proof of concept of the lattice reduction and were able to perform the plaintext recovery using the LLL algorithm in sage[55] with a negligible run time of less than 0.01 seconds (running on a Intel Corei7-4790 CPU @ 3.6GHz).

B. Parallelization of the Bleichenbacher Attack.

The Bleichenbacher attack can also be parallelized in the same way as we have shown for the Manger attack. We assume k parallel attacks. For each attack we start with a different blinding value, such that for attack number j we know that $2B < s_j^0 < 3B - 1$. After i adaptive queries we learn that $a_i < s_j^i < bi$.² Using this information we can recover the plaintext as we have done for the Manger attack.

Analyzing the Parallel Attack. As the Bleichenbacher attack has a much higher query complexity than the Manger attack, we will require a large number of servers to attack. However, if we have k servers, running k attacks in parallel is very inefficient, due to the high cost of the first blinding

phase. Instead we use the fact that each adaptive step of the attack includes many queries that can be done in parallel. We start by using all servers for multiple parallel queries until we find a small number of blinded values (e.g. 5 as in the Manger attack). We then split the k servers evenly between the blinded values to create multiple attacks. For each blinded value, multiple servers will be used to run the parallel queries required for each adaptive step.

VIII. DISCUSSION AND CONCLUSIONS

In this work we have answered negatively the question "Are modern implementations of PKCS #1 v1.5 secure against padding oracle attacks?". The systemic re-discovery of Bleichenbacher's attack on RSA PKCS #1 v1.5 encryption over the last 20 years has shown that the mitigations requirements are unrealistic towards developers. Among the nine popular implementations we surveyed, only two successfully survived our analysis. The insistence that protocols preserve this broken padding standard still have consequences today, reaching even the latest version of TLS 1.3 released in August 2018.

A. Recommendation for mitigations

As we have seen, it is very hard to implement a completely secure and side channel free PKCS #1 v1.5 based RSA key exchange for TLS. We propose several recommendations to help reduce the protocol's vulnerability to our attacks.

Deprecation of RSA key exchange. The safest mitigation is to deprecated the RSA key exchange and switch to (Elliptic Curve) Diffie-Hellman key exchanges. This might be hard due to backward compatibility issues.

Certificate separation. If RSA key exchanges must be supported, it should be done with a dedicated public key that does not allow for signing. If multiple versions of TLS are supported, keys should not be re-used across versions in order to prevent downgrade attacks. If multiple TLS servers are used, each server should use a different public key if possible to prevent parallelization of the attack.

Constant-time code and safe API. The decryption code should be constant-time, with no branching or memory accesses depending on the plaintext (e.g. as achieved in the BoringSSL and BearSSL code). Most of the APIs that we have seen to the decryption function do not allow the calling code to pass the required plaintext size. APIs that pass the expected plaintext length to the unpadding function are safer, since they make it easier to implement the code in constant-time, and when a remaining side channel exists, it will result in a "weaker" Oracle that greatly increases the amount of time required for an attack.

Using large RSA keys. The minimal threshold for decryption using Bleichenbacher and Manger type attacks is $o(\text{number of bits in } N)$ of consecutive calls to the oracle. Using larger keys (at least 2048bit) will require more time for the attack and might make MiTM attack less practical.

Handshake timeouts. It is harder to do a MiTM attack when the TLS handshake timeout is very short. Clients should use

²With low probability we might have more than one possible domain, and in that case we can take the domain from one of the previous queries

short TLS timeouts, and make sure they are resilient to any attack that can lengthen the timeout (such as the TLS warning alerts attack against Firefox[4]).

Speed limitation. As RSA key exchanges are only a small fraction of today's TLS traffic[1, 45], limiting the speed of allowed RSA decryptions makes MiTM attacks less practical.

Dedicated hardware for running sensitive cryptographic code. Side channel attacks are extremely difficult to defend against. Critical and sensitive operations such as private key decryption should not be run on a hardware shared with other code if possible.

B. Future work

Timeouts in TLS client. As we have seen in this work and previous works [4], the possibility of doing some MiTM attacks depends strongly on the amount of time the attacker has before the client gives up on the handshake. Clients that have long handshake timeouts (e.g. curl and git) or are vulnerable to a "timeout extension" attack (e.g. Firefox) put their users at risk. A systematic review of different client's timeouts configuration and their resilience to "timeout extension" attacks is required.

Keyless TLS implementations. Many (often private) TLS implementations segregate private key operations from the protocol implementation by having a keyless server responding to signature and decryption requests from keyless clients. PKCS #1 v1.5 verification is not always done from the keyless server and decrypted ciphertexts of variable-length passed to the keyless clients can be passively observed from a privileged network position. A review of available implementations and standards (like LURK [43]) is needed.

REFERENCES

- [1] "The ICSI Notary," <http://notary.icsi.berkeley.edu/#connection-cipher-details>.
- [2] O. Aciğmez, "Yet another microarchitectural attack: Exploiting I-Cache," in *CSAW*, 2007.
- [3] O. Aciğmez, S. Gueron, and J. Seifert, "New branch prediction vulnerabilities in OpenSSL and necessary software countermeasures," in *IMA Int. Conf.*, 2007.
- [4] D. Adrian, K. Bhargavan, Z. Durumeric, P. Gaudry, M. Green, J. A. Halderman, N. Heninger, D. Springall, E. Thomé, L. Valenta, B. VanderSloot, E. Wustrow, S. Z. Béguelin, and P. Zimmermann, "Imperfect forward secrecy: How Diffie-Hellman fails in practice," in *CCS*, 2015.
- [5] N. J. AlFardan and K. G. Paterson, "Lucky thirteen: Breaking the TLS and DTLS record protocols," in *IEEE SP*, 2013, pp. 526–540.
- [6] T. Allan, B. B. Brumley, K. E. Falkner, J. van de Pol, and Y. Yarom, "Amplifying side channels through performance degradation," in *ACSAC*, 2016.
- [7] R. Bardou, R. Focardi, Y. Kawamoto, L. Simionato, G. Steel, and J. Tsay, "Efficient padding oracle attacks on cryptographic hardware," in *CRYPTO*, 2012.
- [8] M. Ben-Or, B. Chor, and A. Shamir, "On the cryptographic security of single RSA bits," in *STOC*.
- [9] D. J. Bernstein, "Cache-timing attacks on AES," 2005.
- [10] D. Bleichenbacher, "Chosen ciphertext attacks against protocols based on the RSA encryption standard PKCS #1," in *CRYPTO*, 1998.
- [11] H. Böck, J. Somorovsky, and C. Young, "Return of Bleichenbacher's oracle threat (ROBOT)," in *USENIX Security*, 2018.
- [12] F. Brasser, U. Müller, A. Dmitrienko, K. Kostiaainen, S. Capkun, and A. Sadeghi, "Software grand exposure: SGX cache attacks are practical," in *WOOT*, 2017.
- [13] J. V. Bulck, F. Piessens, and R. Strackx, "SGX-Step: A practical attack framework for precise enclave execution control," in *SysTEX@SOSP*, 2017.
- [14] —, "Nemesis: Studying microarchitectural timing leaks in rudimentary CPU interrupt logic," in *CCS*, 2018.
- [15] S. Checkoway, J. Maskiewicz, C. Garman, J. Fried, S. Cohny, M. Green, N. Heninger, R.-P. Weinmann, E. Rescorla, and H. Shacham, "A systematic analysis of the Juniper Dual EC incident," in *CCS*, 2016.
- [16] T. Dierks and C. Allen, "The TLS Protocol Version 1.0," RFC 2246, Jan. 1999.
- [17] T. Dierks and E. Rescorla, "The Transport Layer Security (TLS) Protocol Version 1.1," RFC 4346, Apr. 2006.
- [18] —, "The Transport Layer Security (TLS) Protocol Version 1.2," RFC 5246, Aug. 2008.
- [19] C. Disselkoen, D. Kohlbrenner, L. Porter, and D. M. Tullsen, "Prime+Abort: A timer-free high-precision L3 cache attack using intel TSX," in *USENIX Security*, 2017.
- [20] T. Duong and J. Rizzo, "Here come the \oplus ninjas," 2011.
- [21] D. Evtvushkin, D. Ponomarev, and N. B. Abu-Ghazaleh, "Understanding and mitigating covert channels through branch predictors," *TACO*, vol. 13, no. 1, 2016.
- [22] D. Evtvushkin, D. V. Ponomarev, and N. B. Abu-Ghazaleh, "Jump over ASLR: attacking branch predictors to bypass ASLR," in *MICRO*, 2016.
- [23] D. Evtvushkin, R. Riley, N. B. Abu-Ghazaleh, and D. Ponomarev, "BranchScope: A new side-channel attack on directional branch predictor," in *ASPLOS*, 2018.
- [24] Q. Ge, Y. Yarom, D. Cock, and G. Heiser, "A survey of microarchitectural timing attacks and countermeasures on contemporary hardware," *J. Cryptographic Engineering*, vol. 8, no. 1, 2018.
- [25] Q. Ge, Y. Yarom, and G. Heiser, "No security without time protection: We need a new hardware-software contract," in *APSys*, Aug. 2018.
- [26] D. Genkin, L. Valenta, and Y. Yarom, "May the fourth be with you: A microarchitectural side channel attack on several real-world applications of Curve25519," in *CCS*, 2017.
- [27] D. Genkin, L. Pachmanov, E. Tromer, and Y. Yarom, "Drive-by key-extraction cache attacks from portable code," in *ACNS*, 2018.
- [28] D. Gruss, R. Spreitzer, and S. Mangard, "Cache template attacks: Automating attacks on inclusive last-level

- caches,” in *USENIX Security*, 2015.
- [29] D. Gruss, C. Maurice, K. Wagner, and S. Mangard, “Flush+Flush: A fast and stealthy cache attack,” in *DIMVA*, 2016.
- [30] M. S. Inci, B. Gülmözoglu, G. Irazoqui, T. Eisenbarth, and B. Sunar, “Cache attacks enable bulk key recovery on the cloud,” in *CHES*, 2016.
- [31] Intel, “Speculative execution side channel mitigations,” May 2018.
- [32] G. Irazoqui, T. Eisenbarth, and B. Sunar, “S\$A: A shared cache attack that works across cores and defies VM sandboxing - and its application to AES,” in *IEEE SP*, 2015, pp. 591–604.
- [33] G. Irazoqui, M. S. Inci, T. Eisenbarth, and B. Sunar, “Lucky 13 strikes back,” in *ASIA CCS*, 2015.
- [34] T. Jager, S. Schinzel, and J. Somorovsky, “Bleichenbacher’s attack strikes again: Breaking PKCS#1 v1.5 in XML encryption,” in *ESORICS*, 2012.
- [35] T. Jager, J. Schwenk, and J. Somorovsky, “On the security of TLS 1.3 and QUIC against weaknesses in PKCS#1 v1.5 encryption,” in *CCS*, 2015.
- [36] V. Klíma, O. Pokorný, and T. Rosa, “Attacking RSA-based sessions in SSL/TLS,” in *CHES*, 2003.
- [37] P. Kocher, J. Horn, A. Fogh, D. Genkin, D. Gruss, W. Haas, M. Haburg, M. Lipp, S. Mangard, T. Prescher, M. Schwartz, and Y. Yarom, “Spectre attacks: Exploiting speculative execution,” in *IEEE SP*, 2019.
- [38] S. Lee, M. Shih, P. Gera, T. Kim, H. Kim, and M. Peinado, “Inferring fine-grained control flow inside SGX enclaves with branch shadowing,” in *USENIX Security*, 2017.
- [39] A. K. Lenstra, H. W. Lenstra, and L. Lovász, “Factoring polynomials with rational coefficients,” *Mathematische Annalen*, vol. 261, no. 4, 1982.
- [40] F. Liu, Y. Yarom, Q. Ge, G. Heiser, and R. B. Lee, “Last-level cache side-channel attacks are practical,” in *IEEE SP*, 2015.
- [41] J. Manger, “A chosen ciphertext attack on RSA optimal asymmetric encryption padding (OAEP) as standardized in PKCS #1 v2.0,” in *CRYPTO*, 2001.
- [42] C. Meyer, J. Somorovsky, E. Weiss, J. Schwenk, S. Schinzel, and E. Tews, “Revisiting SSL/TLS implementations: New Bleichenbacher side channels and attacks,” in *USENIX Security*, 2014.
- [43] D. Migault and I. Boureau, “LURK extension version 1 for (D)TLS 1.2 authentication,” IETF, Internet-Draft draft-mgmt-lurk-tls12-01, 2018.
- [44] A. Moghimi, G. Irazoqui, and T. Eisenbarth, “Cachezoom: How SGX amplifies the power of cache attacks,” in *CHES*, 2017.
- [45] Mozilla, “SSL handshake key exchange algorithm for full handshake,” <https://mzl.la/2BQjcMO>.
- [46] OpenSSL, “RSA_public_encrypt,” https://www.openssl.org/docs/man1.0.2/crypto/RSA_private_decrypt.html.
- [47] Y. Oren, V. P. Kemerlis, S. Sethumadhavan, and A. D. Keromytis, “The spy in the sandbox: Practical cache attacks in JavaScript and their implications,” in *CCS*, 2015.
- [48] D. A. Osvik, A. Shamir, and E. Tromer, “Cache attacks and countermeasures: The case of AES,” in *CT-RSA*, 2006.
- [49] C. Percival, “Cache missing for fun and profit,” in *Proceedings of BSDCan*, 2005.
- [50] T. Ristenpart, E. Tromer, H. Shacham, and S. Savage, “Hey, you, get off of my cloud: exploring information leakage in third-party compute clouds,” in *CCS*, 2009.
- [51] R. L. Rivest, A. Shamir, and L. M. Adleman, “A method for obtaining digital signatures and public-key cryptosystems,” *Commun. ACM*, vol. 21, no. 2, 1978.
- [52] E. Ronen, K. G. Paterson, and A. Shamir, “Pseudo constant time implementations of TLS are only pseudo secure,” in *CCS*, 2018.
- [53] *PKCS #1 v2.2: RSA Cryptography Standard*, RSA Laboratories, 2012.
- [54] S. Schmidt, “Introducing s2n, a new open source tls implementation,” 2015.
- [55] The Sage Developers, *SageMath, the Sage Mathematics Software System (Version 8.3)*, www.sagemath.org, 2018.
- [56] Y. Tsunoo, T. Saito, T. Suzaki, M. Shigeri, and H. Miyauchi, “Cryptanalysis of DES implemented on computers with cache,” in *CHES*, 2003.
- [57] Y. Xiao, M. Li, S. Chen, and Y. Zhang, “STACCO: differentially analyzing side-channel traces for detecting SSL/TLS vulnerabilities in secure enclaves,” in *CCS*, 2017.
- [58] M. Yan, C. W. Fletcher, and J. Torrellas, “Cache telepathy: Leveraging shared resource attacks to learn DNN architectures,” *CoRR*, vol. abs/1808.04761, 2018.
- [59] Y. Yarom, “Mastik: A micro-architectural side-channel toolkit,” cs.adelaide.edu.au/~yval/Mastik/Mastik.pdf, 2017.
- [60] Y. Yarom and K. Falkner, “FLUSH+RELOAD: A high resolution, low noise, L3 cache side-channel attack,” in *USENIX Security*, 2014.
- [61] X. Zhang, Y. Xiao, and Y. Zhang, “Return-oriented Flush-Reload side channels on ARM and their implications for Android devices,” in *CCS*, 2016.
- [62] Y. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart, “Cross-tenant side-channel attacks in PaaS clouds,” in *CCS*, 2014.

APPENDIX A VULNERABILITIES DESCRIPTION

A. OpenSSL TLS Implementation

Perhaps aware of the side channel issues in its RSA decryption API, OpenSSL does not use the code described in [Section IV-B](#) for its own TLS implementation. Instead, OpenSSL reimplemented the RSA PKCS #1 v1.5 padding verification as part of its TLS protocol code. This constant time implementation does not appear to be vulnerable to a cache-based padding oracle attack. However, OpenSSL’s code does

```

1 int BN2binpad((bn, to)){
2 //bn is big number (storing the RSA plaintext)
3 //to is the output buffer
4 //BN_BYTES is the number of bytes in each bn word
5
6 i = BN_num_bytes(bn);
7 tolen=i
8 while (i--) {
9     l = bn[i / BN_BYTES];
10    *(to++) = (unsigned char)
11    ( l>> (8 * (i % BN_BYTES))) & 0xff;
12 }
13 return tolen;
14 }

```

Listing 5. Pseudocode of big number serialization functions

contain two side channel vulnerabilities in its data conversion and raw RSA decryption routines which we now describe.

Leaky Data Conversion. As mentioned in Section IV, the big numbers representing the RSA ciphertext and plaintext are typically saved as an array of 32-bit words, while the result of the PKCS #1 v1.5 padding is an array of bytes. To convert the data from one representation to the other, OpenSSL uses a serialization function which takes as input a big number and serializes it into a byte array (where index 0 is the most significant byte). In order not to create a padding oracle, it is important that the serialization function be written in a constant-time manner, and not leak the length of the RSA plaintext during the serialization process.

The pseudocode of OpenSSL’s serialization function is presented in Listing 5. Notice the while loop in Line 8, which performs as many iterations as the number of non-zero bytes of the RSA plaintext, resulting in an extremely efficient Manger-type padding oracle. Traditionally, mounting such precise microarchitectural attacks is difficult, as a single loop iteration takes less time than the channel’s temporal resolution. However, recent works [13, 14, 44] have shown that mounting high precision side channel attacks is possible in the case of trusted execution environments (e.g., Intel SGX), often with cycle-accurate resolution.

B. Amazon s2n

S2n is Amazon’s implementation of the TLS protocol, used as part of Amazon Web Services. It simplifies the OpenSSL TLS implementation, removing uncommon and deprecated TLS configurations. The implementation of RSA decryption (Listing 6) invokes the OpenSSL `RSA_private_decrypt` API function to process and remove the PKCS #1 v1.5 padding (Line 6). We have already discussed the weakness due to the use of the OpenSSL function (Section IV-B). We now discuss another vulnerability in the s2n code.

Leaky PKCS #1 v1.5 Verification. In case the decryption and PKCS #1 v1.5 verification succeeds and the output is of the expected length, s2n copies the data to the output array (Line 7). Moreover, the decision of whether to copy and the copy itself is done in constant time to avoid leaking the result of the result of the PKCS #1 v1.5 unpadding.

However, the s2n API relies on the error status returned from OpenSSL to identify padding failures or mis-formatted

output. Thus, s2n uses an if macro, which compiles to a conditional branch (see Line 8), which yields an FFFT oracle.

C. MbedTLS

MbedTLS aims at providing a portable, easy to use and to read implementation of the TLS protocol and is designed primarily to be used in low powered embedded devices. We have identified vulnerabilities in both the data conversion and the PKCS #1 v1.5 verification stages of the mbedTLS implementation which we now describe.

Leaky PKCS #1 v1.5 Verification. Listing 7 shows the relevant parts of the mbedTLS PKCS #1 v1.5 verification. For brevity we omit the padding format and plaintext length validation, which execute in constant-time. The rest of the code, however uses conditional branches to handle padding validation failures (Lines 7–10) and incorrect plaintext length (Lines 12–15). Thus, despite the constant-time validation, the following form of oracles are still exposed.

- **Potentially Leaky Comparison.** First, the comparison in Line 6 may be implemented using conditional statements, which would leak via branch prediction. This does not happen in our test environment, where the comparison is implemented using a conditional `set` instruction, which to the best of our knowledge executes in constant-time. However without a guarantee that the compiler will use a constant-time implementation there is a potential for a leak in other environments.
- **Length Dependant Branches.** Both if statements in Lines 7 and 12 can be exploited for a branch prediction attack. The former allows a FFFT Bleichenbacher oracle and the latter allows an FFFT oracle variant. In fact, the oracle is slightly stronger than a standard FFFT oracle because the test is one sided, i.e. it only checks for maximum size instead of checking for exact size.
- **Length Dependant Early Termination.** Finally, due to early termination on bad inputs, the code that copies to the output (Line 18) is only executed if the plaintext is PKCS #1 v1.5 conforming. Thus we can implement an FFFT oracle via an instruction cache attack, monitoring either the call to `memcpy` or the code of `memcpy` itself.

Leaky Data Conversion. The last step in the implementation of RSA decryption in mbedTLS is to copy the plaintext to the output. As discussed in Section IV, there is no a-priory method for determining the plaintext’s length, and applications can only determine the length after decryption. To determine the length, mbedTLS scans the words that represent the plaintext from the most significant to the least significant, looking for a non-zero word. In a padding oracle attack, this is very likely to be the first word of the plaintext. MbedTLS then scans the bits of the word to find the most significant non-zero bit. This scan, shown in Listing 8, loops over the bits, from the most significant to the least significant (Line 7), checking for a non-zero bit (Line 8). An adversary that can count the number of iterations executed can learn the leading number of zero bits, which can be used for a Manger type oracle. As in

```

1 int s2n_rsa_decrypt(priv, in, out){
2   unsigned char intermediate[4096];
3   const s2n_rsa_private_key *key = &priv->key.rsa_key;
4   S2N_ERROR_IF(s2n_rsa_private_encrypted_size(key) > sizeof(intermediate), S2N_ERR_NOMEM);
5   S2N_ERROR_IF(out->size > sizeof(intermediate), S2N_ERR_NOMEM);
6   int r = RSA_private_decrypt(in->size, in->data, intermediate, key->rsa, RSA_PKCS1_PADDING);
7   GUARD(s2n_constant_time_copy_or_dont(out->data, intermediate, out->size, r != out->size));
8   S2N_ERROR_IF(r != out->size, S2N_ERR_SIZE_MISMATCH);
9   return 0;
10 }

```

Listing 6. Pseudocode of Amazon s2n’s wrap for OpenSSL’s API

```

1 int mbedtls_rsa_rsaes_pkcs1_v15_decrypt(
2   ilen, olen, input, output, output_max_len) {
3   ...
4   //Omitted code checks for valid padding and
   //length of decrypted plaintext
5
6   bad |= ( pad_count < 8 );
7   if( bad ){
8     ret = MBEDTLS_ERR_RSA_INVALID_PADDING;
9     goto cleanup;
10  }
11
12  if( ilen - ( p - buf ) > output_max_len ){
13    ret = MBEDTLS_ERR_RSA_OUTPUT_TOO_LARGE;
14    goto cleanup;
15  }
16
17  *olen = ilen - ( p - buf);
18  memcpy( output, p, *olen );
19  ret = 0;
20
21  cleanup:
22  mbedtls_zeroize( buf, sizeof( buf ) );
23  return( ret );
24 }

```

Listing 7. MbedTLS’s unpadding function

```

1 size_t mbedtls_clz( x ){
2   // x is the RSA decrypted plaintext
3   // biL is the number of bits in limb (typ. 64)
4   size_t j;
5   mask = 1 << (biL - 1);
6
7   for( j = 0; j < biL; j++ ){
8     if( x & mask ) break;
9     mask >>= 1;
10  }
11  return j;
12 }

```

Listing 8. MbedTLS’s bit length checking function function

Appendix A-A, such attacks are unfeasible for unprivileged adversaries, but can be performed by a root adversary attacking a code running in trusted execution environment (e.g., Intel SGX). Finally, we note that the adversary only needs to determine whether the loop body gets executed for implementing an Interval oracle (see [Section II-E](#)).

D. Mozilla NSS

Mozilla’s Network Security Services (NSS) library is the cryptographic engine often used in applications developed by the Mozilla project. NSS implements countermeasures for padding oracle attacks, however, the TLS code ignores the possibility of leakage through microarchitectural channels

```

1 mp_to_fixlen_octets(mp, str, length)
2 {
3   // mp is a number encoded in little endian
4   // str is an array of length bytes containing
5   // a big endian encoding of mp
6   int ix, pos = 0;
7   unsigned int bytes;
8   bytes = mp_unsigned_octet_size(mp);
9   /* place any needed leading zeros */
10  for (; length > bytes; --length) {
11    *str++ = 0;
12  }
13  .../* code for convering a little-endian large
14     * number mp into a big-endian fixed-length
15     * byte array str (omitted for brevity) */
16 }

```

Listing 9. Data Conversion in NSS

Consequently, the TLS implementation exposes padding oracle in each of the three stages of handling PKCS #1 v1.5 padding.

Leaky Data Conversion. [Listing 9](#) shows a leak in the data conversion stage. The code is the start of the function `mp_to_fixlen_octets`, which converts a large number into a fixed-length byte array. The function first determines the number of bytes required for storing the number (Line 8). Next, it zero-pads the output byte array, so that the final output is exactly `length` bytes (Lines 10–12). Finally, it converts the large number `m` from its little-endian representation to a big-endian byte array representation (omitted for brevity).

Unfortunately, `mp_to_fixlen_octets` does not perform the padding in constant time, thus leaking the number of leading zeros in the RSA decrypted plaintext to an adversary that can count (via the cache side channel) the number of iterations in the loop in Lines 10–12. Furthermore, a branch prediction attack can determine whether the body of the loop executed, allowing a Manger-type oracle.

Leaky PKCS #1 v1.5 Verification. We now describe the leaks from the PKCS #1 v1.5 verification code in NSS ([Listing 10](#)). The code performs a textbook verification of the PKCS #1 v1.5 format, e.g. Lines 10 and 11 check the values of the first two bytes in the message.

Unfortunately, the code in [Listing 10](#) terminates early in case of verification failure. Thus, using a branch prediction attack to monitor any of the if statements in the code yields an TTTT-type padding oracle. Moreover, in case that the checks in Lines 10 and 11 are compiled into two different branches this can allow for a Manger type Oracle. Furthermore, as in [Appendix A-C](#), monitoring the call to `PORT_Memcpy` (Line 24) using a cache side channel yields a stronger variant

```

1 RSA_DecryptBlock(key, output, outputLen,
2                 maxOutputLen, input, inputLen)
3 {
4     ...
5     rv = RSA_PrivateKeyOp(key, buffer, input);
6     if (rv != SECSuccess)
7         goto loser;
8
9     /* XXX(rsleeve): Constant time */
10    if (buffer[0] != RSA_BLOCK_FIRST_OCTET ||
11        buffer[1] != RSA_BlockPublic) {
12        goto loser;
13    }
14    *outputLen = 0;
15    for (i = 2; i < modulusLen; i++) {
16        if (buffer[i] == RSA_BLOCK_AFTER_PAD_OCTET) {
17            *outputLen = modulusLen - i - 1;
18            break;
19        }
20    }
21    if (*outputLen == 0)
22        goto loser;
23    ...
24    PORT_Memcpy(output, buffer + modulusLen - *
25                outputLen, *outputLen);
26    return SECSuccess;
27 loser:
28    PORT_Free(buffer);
29 failure:
30    return SECFailure;
31 }

```

Listing 10. NSS’s PKCS #1 v1.5 Verification function

```

1 wc_RsaFunctionSync(in, inLen, out, outLen, key)
2 {
3     ... // code for performing RSA decryption of in
4     // result is stored in temp
5     if (ret == 0) {
6         len = mp_unsigned_bin_size(tmp);
7         while (len < keyLen) {
8             *out++ = 0x00;
9             len++;
10        }
11        ...
12    }
13    ...
14 }

```

Listing 11. WolfSSL’s RSA decryption conversion

of FTTT-type padding oracle, as it only checks for zero anywhere after the first 2 bytes.

Leaky Padding Oracle Mitigations. Finally, as in OpenSSL (Listing 2), the NSS code responsible for mitigating padding oracle attacks checks the results of the PKCS #1 v1.5 verification procedure using an if statement that translates to a conditional branch. Thus, monitoring this branch as done for Section IV-B results in a FFFF-type padding oracle.

E. WolfSSL

WolfSSL is a TLS library aimed at embedded devices. As in NSS, the WolfSSL code exposes oracles in all stages of PKCS #1 v1.5 handling.

Leaky RSA Decryption Routine. After performing RSA decryption, WolfSSL pads the plaintext to the length of the RSA modulus (Lines 7–10 in Listing 11) using a while loop.

```

1 void nettle_mpz_to_octets(length, *s, x, sign) {
2     // convert x in little endian big number to
3     // a big endian byte array representation s
4     // of length bytes
5     uint8_t *dst = s + length - 1;
6     size_t size = mpz_size(x);
7     size_t i;
8
9     for (i = 0; i < size; i++) {
10        mp_limb_t limb = mpz_getlimbn(x, i);
11        size_t j;
12        for (j = 0; length && j < sizeof(mp_limb_t); j++) {
13            *dst-- = sign ^ (limb & 0xff);
14            limb >>= 8;
15            length--;
16        }
17    }
18    if (length) memset(s, sign, length);
19 }

```

Listing 12. GnuTLS’s Data Conversion function

The number of iterations this loop performs leaks the number of leading zero bytes, exposing a Manger oracle.

Leaky PKCS #1 v1.5 Verification and Padding Oracle Mitigations. Additionally WolfSSL uses a naive, variable time code for PKCS #1 v1.5 verification, leaking Manger and FTTT-type padding oracles. Moreover, the padding oracle mitigation code leaks FTTT- and FFFF-type padding oracles through the microarchitectural channels.

F. GnuTLS

GnuTLS is another popular implementation of the TLS protocol. Like WolfSSL and NSS, GnuTLS does not use constant time code for the PKCS #1 v1.5 verification, resulting in numerous side-channel-observable padding oracles.

Leaky Data Conversion. To convert RSA-decrypted plaintext from a little-endian big number format to big-endian byte array format, GnuTLS uses code from the Nettle cryptographic library³. Listing 12 shows the data conversion code in Nettle. Line 18 conditionally calls `memset` when there are leading zeros in the plaintext, exposing a Manger oracle.

Leaky PKCS #1 v1.5 Verification. GnuTLS also relies on leaky Nettle for PKCS #1 v1.5 verification (Listing 13). The branch in Line 10 allows for a Manger type oracle or a TTTT oracle. The conditional call to `memcpy` in Line 15 exposes an FTTT oracle.

Leaky Padding Oracle Mitigations. The GnuTLS padding oracle mitigation code is also not constant-time, see Listing 14 for a simplified version. In particular, the branches in Lines 7 and 12 yield a FTTT Bleichenbacher oracle. Another potential security issue in the code present in Listing 14 is the comment “we do not need strong random numbers here” (Line 3). We note that predicting the random session key used for padding oracle mitigation, renders the mitigation ineffective. This is since the attacker can use this session key to generate the correct client finish message, thereby causing the server to complete the TLS handshake. This results in

³<https://www.lysator.liu.se/~nisse/nettle/>

```

1 int pkcs1_decrypt(key_size, m, length, message) {
2     TMP_GMP_DECL(em, uint8_t);
3     uint8_t *terminator;
4     size_t padding;
5     size_t message_length;
6     int ret;
7     TMP_GMP_ALLOC(em, key_size);
8     nettle_mpz_get_str_256(key_size, em, m);
9     /* Check format */
10    if (em[0] || em[1] != 2) {
11        ret = 0;
12        goto cleanup;
13    }
14    ...
15    memcpy(message, terminator+1, message_length);
16    *length = message_length;
17    ret = 1;
18    cleanup:
19    TMP_GMP_FREE(em);
20    return ret;
21 }

```

Listing 13. GnuTLS's PKCS #1 v1.5 verification

```

1 int proc_rsa_client_kx(session, data) {
2     ...
3     // we do not need strong random numbers here.
4     ret = gnutls_rnd(GNUTLS_RND_NONCE, rndkey.data,
5                     rndkey.size);
6     ...
7     ret = gnutls_privkey_decrypt_data(session->
8         internals.selected_key, 0, &data, &plaintext);
9     if (ret < 0 || plaintext.size != GNUTLS_MASTER_SIZE) {
10        randomize_key = 1;
11        ...
12    }
13    ...
14    if (randomize_key != 0) {
15        session->key.key.data = rndkey.data;
16        session->key.key.size = rndkey.size;
17        rndkey.data = NULL;
18    } else {
19        session->key.key.data = plaintext.data;
20        session->key.key.size = plaintext.size;
21    }
22    return ret;
23 }

```

Listing 14. Pseudocode of GnuTLS's padding oracle mitigation

a remote Bleichenbacher FFTT oracle that does not require any side channel leakage. We leave further exploration of bad randomness used in padding oracle mitigations to future work.