

A TECHNICAL NOTE ON UNDERSTANDING FIPS 204 (ML-DSA)

Bo Lin



Version history

Version	Date	Author	Status
A	02 Feb 2026	B LIN	Release

Change history

Version	Changes

BO

Table of Contents

1. Overview and notations	4
2. Learning with Errors (LWE)	4
3. Two simplified FIPS 204 versions	5
3.1. Key generation – KeyGen()	5
3.2. Signing – Sign_internal()	6
3.3. Verifying – Verify_internal()	7
4. The refinement in FIPS 204	7
5. Summary	8

BOUNN

1. Overview and notations

An initial challenge for those who directly dive into the FIPS 204 ML-DSA specification (*FIPS 204, Module-Lattice-Based Digital Signature Standard, August 13, 2024, <https://csrc.nist.gov/pubs/fips/204/final>*) is that data compression and security checks are interleaved in various functions in the lattice-based signature scheme, making it difficult to map a conceptual lattice-based signature scheme into the FIPS 204 ML-DSA specification. The motivation of this technical note is to bridge this gap to overcome the challenge.

This technical note begins with a brief description of the hard problem in lattice-based cryptography by examining the difference between the complexity of solving $Ax = t \bmod p$ and the complexity of solving $Ax + e = t \bmod p$ (a.k.a. the Learning with Errors (LWE) problem), where p is a pre-defined prime, A is a known $k \times l$ matrix, t a known k -dimensional vector, x an unknown l -dimensional vector, and e an unknown l -dimensional vector.

Then, it proceeds by tabulating two simplified versions of the FIPS 204 ML-DSA – one is a simple math version (SMV) and the other a simplified specification version (SSV) – to explain the FIPS 204 ML-DSA specification. The SMV intends to explain the math fundamentals on how a lattice-based digital signature works while the SSV serves a stepping-stone from the SMV to the FIPS 204 ML-DSA specification.

The technical note ends with a simplified flow-chart for Algorithm 7 ML-DSA.Sign_internal() of the FIPS 204 ML-DSA specification to visualise some key steps.

In this technical note, a bold letter, such as A or x , stands for a matrix or a vector. Its entry, as that in the FIPS 204 ML-DSA specification, is a polynomial in $Z_p[X]/(X^{256} + 1)$. On the other hand, a regular letter, such as c , stands for a single element. It can be a polynomial in $Z_p[X]/(X^{256} + 1)$, a data string, or an integer.

If a symbol is the same as that in the FIPS 204 ML-DSA specification, the symbol follows its definition in the specification to make the cross-reference easy.

2. Learning with Errors (LWE)

For $Ax = t \bmod p$, the system of equations can be solved, or the x can be *learned*, efficiently by solving $(A - t) \cdot \begin{pmatrix} x \\ 1 \end{pmatrix} = \mathbf{0}$. However, if the $Ax = t \bmod p$ is changed to $Ax + e = t \bmod p$ by adding a k -dimensional vector e whose components are small random numbers, i.e., by adding small random errors in the system, this system of equations can be written as $(A I_k - t) \cdot \begin{pmatrix} x \\ e \\ 1 \end{pmatrix} = \mathbf{0} \bmod p$, where the I_k is an identity matrix of order k .

The “Given A and t to find x and e ” becomes an instance of Learning with Errors (LWE). It can be shown that the vector $(x \ e \ 1)$ is a shortest vector in the lattice pertaining to the A , meaning that solving the LWE problem yields to solving the Shortest Vector Problem (SVP) which is extremely difficult for general cases.

There are several variants of LWE. NIST selected CRYSTALS-Dilithium (<https://pq-crystals.org/dilithium/>), a digital signature scheme based on module LWE, to standardise the Post-Quantum Cryptography digital signature algorithm, ML-DSA (Module-Lattice-Based Digital Signature Standard).

In the FIPS 204 ML-DSA specification, the x is denoted as s_1 , e as s_2 . As a result, the $Ax + e = t \bmod p$ is denoted as $As_1 + s_2 = t \bmod q$ with $q = 8380417$.

3. Two simplified FIPS 204 versions

As mentioned in section 1. *Overview and notations*, the SMV and SSV are tabulated for contrast. The SMV provides a mathematics proof for the correctness of the signature scheme and the SSV serves as a stepping-stone to the FIPS 204 specification.

3.1. Key generation – KeyGen()

#	SMV (Simple Math Version)	SSV (Simplified Spec Version)
1	Generate $A \in R_q^{k \times l}$. That is, the A is a $k \times l$ matrix with each entry is a polynomial in $Z_q[X]/(X^{256} + 1)$. Simply put, an entry of A is a $GF(q)$ polynomial modulo $X^{256} + 1$ with its coefficients being random numbers in $[0, q - 1]$. Note, in the FIPS 204 ML-DSA specification, \widehat{A} , the A 's “spectrum” representation, is generated to facilitate the NTT (Number Theory Transform) operation for polynomial multiplications because a random A 's NTT result, $\widehat{A} = \text{NTT}(A)$, is a random matrix anyway, so the random \widehat{A} can be regarded as $\text{NTT}(A)$ for some random A .	
2	Sample $s_1 \in R_q^l$ with coefficients in $[-\eta, \eta]$ where $\eta = 2$ for ML-DSA-44 while $\eta = 4$ for ML-DSA-65 and ML-DSA-87. That is, the s_1 is a vector of l components and each component is a 255-degree polynomial with coefficients in $[-\eta, \eta]$.	
3	Sample $s_2 \in R_q^k$ with coefficients in $[-\eta, \eta]$ where $\eta = 2$ for ML-DSA-44 while $\eta = 4$ for ML-DSA-65 and ML-DSA-87. The s_2 is a vector of k components and each component is a 255-degree polynomial with coefficients in $[-\eta, \eta]$.	
4	Calculate $t = As_1 + s_2$. Note, the As_1 operation is implemented as $\text{NTT}^{-1}(\widehat{A} \circ \text{NTT}(s_1))$ where the \circ is component-wise multiplication.	
5	–	$t = t_1 \cdot 2^d + t_0 \bmod q$ $t_0 = (t \bmod q) \bmod^{\pm} 2^d$ $d = 13$
6	$pk = (A, t)$, $sk = s_1$	$pk = (A, t_1)$, $sk = (s_1, s_2, t_0)$

It is noted that the s_2 is not part of sk in the SMV, but it is part of sk in the SSV. In addition, the t is split into t_1 and t_0 in the SSV, where the t_1 (higher part of t) is used as part of pk while the t_0 (lower part of t) as part of sk .

In the FIPS 204 ML-DSA specification, the s_2 is used for security check and it is also used with t_0 to make a “hint” vector h in `Sign_internal()`. The vector h later used in `Verify_internal()` to adjust the result from `HighBits()`.

3.2. Signing – Sign_Internal()

#	SMV (Simple Math Version)	SSV (Simplified Spec Version)
1	Sample $\mathbf{y} \in R_q^l$ with small coefficients in $[-\gamma_1 + 1, \gamma_1]$ where $\gamma_1 = 2^{17}$ for ML-DSA-44 while $\gamma_1 = 2^{19}$ for ML-DSA-65 and ML-DSA-87. That is, the \mathbf{y} is a vector of l components and each component is a 255-degree polynomial with coefficients in $[-\gamma_1 + 1, \gamma_1]$.	
2	Calculate $\mathbf{w} = \mathbf{A}\mathbf{y}$. Note, the $\mathbf{A}\mathbf{y}$ operation is implemented as $\text{NTT}^{-1}(\hat{\mathbf{A}} \circ \text{NTT}(\mathbf{y}))$ where the \circ is component-wise multiplication.	
3	Calculate $\mathbf{w}_1 = \text{HighBits}(\mathbf{w})$ // the HighBits() is a rounding function	
4	Calculate $\tilde{c} = \text{SHAKE256}(\mu \mathbf{w}_1)$ // the \tilde{c} is a bit-string	
5	Transform the bit-string \tilde{c} to a polynomial c .	
6	Calculate $\mathbf{z} = \mathbf{y} + cs_1$ Note, the cs_1 operation is implemented as $\text{NTT}^{-1}(\text{NTT}(c) \circ \text{NTT}(s_1))$ where the \circ is component-wise multiplication.	
7	–	Use s_2 and γ_1 to check security. If fails, go to #1. Note, the γ_1 is defined in the FIPS 204 ML-DSA specification
8	–	Use t_0 and s_2 to MakeHint() \mathbf{h} . If the \mathbf{h} is not legitimate, go to #1.
9	Output signature $\sigma = (\mathbf{z}, c)$	Output signature $\sigma = \text{sigEncode}(\mathbf{z}, \tilde{c}, \mathbf{h})$

As can be seen in Sign_Internal().SMV, the s_2 is not used to work out σ explicitly, neither in the Sign_Internal().SSV, but it is used for security check and the \mathbf{h} calculation in Sign_Internal().SSV.

Line #4 of Sign_Internal() needs attention that the bit-string \tilde{c} is not suitable for multiplying a polynomial vector s_1 . It needs to be transformed to be a polynomial c as that in line #5. This is implemented in Algorithm 29 SampleBall() of the FIPS 204 ML-DSA specification.

The output of Sign_Internal().SMV $\sigma = (\mathbf{z}, c)$ is conceptual while the output $\sigma = \text{sigEncode}(\mathbf{z}, \tilde{c}, \mathbf{h})$ of Sign_Internal().SSV is a bit-string by compressing and coding \mathbf{z} , \tilde{c} and \mathbf{h} .

3.3. Verifying – Verify_internal()

#	SMV (Simple Math Version)	SSV (Simplified Spec Version)
1	Recover c from \tilde{c} . // \tilde{c} is a bit string and c is a polynomial	
2	–	Calculate $\mathbf{w}'_{\text{Approx}} = \mathbf{A}\mathbf{z} - c\mathbf{t}_1 \cdot 2^d$ Note, the $\mathbf{A}\mathbf{z} - c\mathbf{t}_1 \cdot 2^d$ operation is implemented as $\text{NTT}^{-1}(\hat{\mathbf{A}} \circ \text{NTT}(\mathbf{z}) - \text{NTT}(c) \circ \text{NTT}(\mathbf{t}_1 \cdot 2^d))$ where the \circ is component-wise multiplication.
3	Calculate $\mathbf{w}'_1 = \text{HighBits}(\mathbf{A}\mathbf{z} - c\mathbf{t})$	$\mathbf{w}'_1 = \text{UseHint}(\mathbf{h}, \mathbf{w}'_{\text{Approx}})$
4	Calculate $\tilde{c}' = \text{SHAKE256}(\mu \mathbf{w}'_1)$ // the \tilde{c}' is a bit-string	
5	Output: valid if $\tilde{c} == \tilde{c}'$, else invalid // the \tilde{c} is received and the \tilde{c}' calculated.	

In line #3 of Verify_internal().SMV,

$$\begin{aligned} \mathbf{w}'_1 &= \text{HighBits}(\mathbf{A}\mathbf{z} - c\mathbf{t}) = \text{HighBits}(\mathbf{A}(\mathbf{y} + cs_1) - c(\mathbf{A}s_1 + s_2)) \\ &= \text{HighBits}(\mathbf{A}\mathbf{y} + Acs_1 - Acs_1 - cs_2) \\ &= \text{HighBits}(\mathbf{A}\mathbf{y} - cs_2) \end{aligned}$$

On the other hand, in line #3 of Sign_internal().SMV, $\mathbf{w}_1 = \text{HighBits}(\mathbf{A}\mathbf{y})$. So, when s_2 is small, the rounded result leads to $\mathbf{w}'_1 = \mathbf{w}_1$, resulting in $\tilde{c} = \tilde{c}'$ if the message μ is not tampered.

Line #2 and line #3 of Verify_internal().SSV implement the rounding operation. The higher part of t (i.e., \mathbf{t}_1) is used in line #2 to work out an approximate rounded value of $(\mathbf{A}\mathbf{z} - ct)$ as $\mathbf{w}'_{\text{Approx}}$. The $\mathbf{w}'_{\text{Approx}}$ is further rendered with \mathbf{h} to get \mathbf{w}'_1 in line #3 to achieve the final rounded value of $(\mathbf{A}\mathbf{z} - ct)$. The \mathbf{w}'_1 in line #3 of Verify_internal().SSV matches the \mathbf{w}_1 in line #3 of Sign_internal().

The key point is that although the two results, $\mathbf{A}\mathbf{y}$ and $\mathbf{A}\mathbf{z} - ct$, are different, after the rounding operation, their rounded results, $\text{HighBits}(\mathbf{A}\mathbf{y})$ and $\text{HighBits}(\mathbf{A}\mathbf{z} - ct)$, are the same.

4. The refinement in FIPS 204

By contrasting the SMV and the SSV as illustrated in section

3. Two simplified FIPS 204 versions, the SSV performs the rounding operation by using the higher part of t (i.e., the \mathbf{t}_1) in line #2 of Verify_internal() and then improving the result with the hint vector \mathbf{h} in line #3 of Verify_internal(). The hint vector \mathbf{h} is calculated with the lower part of t (i.e., the \mathbf{t}_0) in line #8 of Sign_internal() which “hints” whether a carry has happened during the signing process. So, the \mathbf{h} is included in the signature to “inform” Verify_internal() to improve the rounded value as shown in line #3 of Verify_internal().

The use of t_1 makes the public key being compressed from $pk = (\mathbf{A}, \mathbf{t})$ to $pk = (\mathbf{A}, \mathbf{t}_1)$ but the \mathbf{h} makes the signature being expanded from $= (\mathbf{z}, \tilde{\mathbf{c}})$ to $\sigma = \text{sigEncode}(\mathbf{z}, \tilde{\mathbf{c}}, \mathbf{h})$. Because an entry of \mathbf{h} only has two possible values: 0 or 1, the \mathbf{h} can be compressed and coded in a small bit-string according to Algorithm 39 MakeHint() and Algorithm 26 sigEncode() in the FIPS 204 ML-DSA specification.

In the FIPS 204 standard, the \mathbf{A} is compressed by having its NTT representation (i.e., $\widehat{\mathbf{A}}$) generated from a public random 32-byte seed ρ . The 32-byte ρ represents an $k \times l$ matrix $\widehat{\mathbf{A}}$ through Algorithm 32 ExpandA() in the FIPS 204 ML-DSA specification. This leads to $pk = (\rho, \mathbf{t}_1)$ that is a significant reduction of the public key size. The reduction on \mathbf{A} and \mathbf{t} much overweighs the small expansion of the signature size by appending the coded \mathbf{h} .

In addition to the data compression above, security check is also considered. The result on line #6 of Sign_internal() is checked with s_2 and γ_1 on line #7. If the signature result does not meet a set of security criteria, the signature is discarded and a new signature is generated until a legitimate signature is generated (see line #23 of Algorithm 7 ML-DSA.Sign_internal() of the FIPS 204 ML-DSA specification).

In the FIPS 204 ML-DSA specification, some loops are not bounded, these unbounded loops seem to cause an infinitive loop in the signature generation, leading to a potential “signature failure”, that is, no signature will be generated, but, as discussed in *Appendix C – Loop Bounds* of the FIPS 204 ML-DSA specification, the failure probability is negligible ($< 2^{-256}$).

5. Summary

This technical note concludes with a simplified flow-chat that illustrates the key steps of the Sign_internal() in the FIPS 204 ML-DSA specification.

