# Castryck-Decru attack on SIKE SIDH (a very rough overview)

Andrew Sutherland

August 1, 2022

## So what exactly happened?

Over the weekend Wouter Castryck and Thomas Decru posted An efficient key recovery attack on SIDH to the Cryptology ePrint archive which describes an attack on the SIKE SIDH protocol that recently advanced to the fourth round of NIST's ongoing Post-Quantum Cryptography (PQC) standardization process.

This is not a theoretical attack. The authors provide Magma code that:

- solves Microsoft's USD 50,000 $IKEp217 challenge in less than five minutes on a single core (as I will demonstrate).
- breaks SIKEp434 (aimed at NIST quantum security level 1) in less than an hour.
- breaks SIKEp503 (level 2), SIKEp610 (level 3), and SIKEp751 (level 5) in about 2, 8, and 20 hours (respectively).

Modulo needing to factor fixed publicly known integers whose bit-size is about half the size of $p$ (where the public keys are elements of $\mathbb{F}_{p^2}$), the algorithm heuristically runs in deterministic polynomial-time on a classical computer.

## So what does this mean?

The SIKE SIDH protocol proposed to NIST is clearly no longer viable.

Beyond that it is surely too soon to say. To quote the excellent summary Steven Galbraith posted on his Elliptic Curve Cryptography blog:

*The correct response to this is not to attempt to minimise the impact, nor to reflexively declare the subject dead. Instead, we should keep our minds open and let the mathematicians work out the implications, wherever they lead.*

My goal tonight is to give a high-level overview of the main ideas and to demonstrate the attack, not to speculate on the future of isogeny-based cyptography.

From the perspective of number theorists seeking to be informed by computation, this is an exciting and welcome result!

# SIKE SIDH Setup

Let us begin by recalling the supersingular isogeny Diffie-Hellman protocol (SIDH), as described in the Supersingular Isogeny Key Encapsulation (SIKE) proposal to NIST.

Fixed public parameters:

- A prime $p = 2^a 3^b - 1$ with $2^a \approx 3^b$
- Supersingular $E_{\text{start}} : y^2 = x^3 + 6x^2 + x$ over $\mathbb{F}_{p^2}$ with $E(\mathbb{F}_{p^2}) = E[2^a 3^b]$.

Notice that $\mathbb{Z}[2i] \subseteq \text{End}(E)$. The attack exploits but it is not essential, any small non-integer endomorphism would work, and it is not clear that even this is necessary. See section 8.3 of the paper for further discussion.

Alice's and Bob's public parameters:

- Alice picks and publishes generators $P_A, Q_A$ for $E[2^a] = \langle P_A, Q_A \rangle$.
- Bob picks and publishes generators $P_B, Q_B$ for $E[3^b] = \langle P_B, Q_B \rangle$.

These play a crucial role in the attack (and have always been a concern for SIDH).

# Diffie-Hellman key exachange

To agree on a shared secret Alice and Bob execute the following protocol:

- Alice generates random $\mathsf{sk}_A \in [0, 2^a)$ and sends $(E_A, \phi_A(P_B), \phi_A(Q_B))$ to Bob, where $\phi_A \colon E_{\mathrm{start}} \to E_A := E/\langle P_A + [\mathsf{sk}_A]Q_A\rangle$.
- Bob generates random $\mathsf{sk}_B \in [0, 3^b)$ and sends $(E_B, \phi_B(P_A), \phi_B(Q_A))$ to Alice, where $\phi_B \colon E_{\mathrm{start}} \to E_B := E/\langle P_B + [\mathsf{sk}_B]Q_B\rangle$.
- Alice computes $j(E_{AB})$ where $E_{AB} := E_B/\langle \phi_B(P_A) + [\mathsf{sk}_A]\phi_B(Q_A)\rangle$.
- Bob computes $j(E_{BA})$ where $E_{BA} := E_A/\langle \phi_A(P_B) + [\mathsf{sk}_B]\phi_A(Q_B)\rangle$.

Now Alice and Bob both know the shared secret $j(E_{AB}) = j(E_{BA})$.

Note that $\phi_A(P_B)$ and $\phi_A(Q_B)$ are sent in the clear, so to learn the shared secret it is enough to compute $\mathsf{sk}_B$, which is what the attack does. It can be adapted to compute $\mathsf{sk}_A$ as well, but this is not necessary: once you know $\mathrm{sk}_B$ you can compute the shared secret the same way Bob does.

## What is the main idea behind the attack?

The attack exploits the fact that if you start from a product of supersingular elliptic curves $C \times E$ and walk the $(\ell, \ell)$-isogeny graph of principally polarized superspecial abelian surfaces, almost all the vertices you encounter will be Jacobians of genus 2 curves ($\approx p^3/2880$ of them); products of elliptic curves are rare ($\approx p^2/288$ of them). If $p$ is cryptographically large you will almost never see another product of elliptic curves once you depart from $C \times E$ on a walk of polynomial length.

But if your destination $A$ happens to be the codomain of an $(\ell^a, \ell^a)$-isogeny whose kernel is a maximal isotropic subgroup of $(C \times E)[\ell^a]$ arising from an anti-isometry $\psi \colon C[\ell^a] \to E[\ell^a]$ associated to an isogeny diamond of order $\ell^a$, then $A$ will be a product of elliptic curves. And at least heuristically, the converse appears to almost always hold (the attack relies on this heuristic).

One can turn this into an effective test that makes it possible to compute Bob's secret key by iteratively guessing successive ternary digits and testing each guess. When $\ell = 2$ the edges we walk are almost all Richelot isogenies, which makes things particularly simple, explicit, and fast, but this is not essential to the attack, as noted in Section 8.1.

# The main workhorse of the algorithm

In their paper Castryck and Decru describe an algorithm that takes inputs:

(i) A prime $p = 2^a 3^b f - 1$ with $2^a \approx 3^b$ (we will assume $2^a > 3^b$ and $f = 1$).

(ii) A supersingular elliptic curve $E_0/\mathbb{F}_{p^2}$ with $E_0(\mathbb{F}_{p^2}) = E_0[2^a 3^b]$ of order $(p+1)^2$.

(iii) Generators $P_0, Q_0$ for $E[2^a] = \langle P_0, Q_0 \rangle$.

(iv) A known cyclic $3^\beta$-isogeny $\tau \colon E_0 \to E_{\text{start}} := y^2 = x^3 + 6x^2 + x$ with $0 \le \beta < b$.

(v) The codomain $E/\mathbb{F}_{p^2}$ of a secret cyclic $3^b$-isogeny $\varphi \colon E_0 \to E$.

(vi) The points $P = \varphi(P_0)$ and $Q = \varphi(Q_0)$ generating $E[2^a] = \langle P, Q \rangle$.

The algorithm outputs:

- A generator of the kernel of $\varphi$ specified as a linear combination of $P_0$ and $Q_0$ (of the form $P_0 + \kappa Q_0$ where $\kappa$ is an integer we compute in base 3).

To determine $\mathsf{sk}_B$ this algorithm is applied iteratively, starting with $\beta = 0$ and $E_0 = E_{\text{start}}$. Each output $\kappa$ gives another few ternary digits (typically one) of $\mathsf{sk}_B$.

## A decision problem

**(1)** Given $E$, $P$, $Q$, does there exist a cyclic $3^b$-isogeny $\varphi\colon E_0 \to E$ such that $\varphi(P_0) = P$ and $\varphi(Q_0) = Q$, with $E_0, P_0, Q_0$ as in (i),(ii),(iii)?

Suppose (1) holds.

Let $c := 2^a - 3^b > 0$ and $x \in [0, 2^a)$ be the inverse of $3^b \bmod 2^a$ so $-xc \equiv 1 \bmod 2^a$.
Let $\gamma\colon E_0 \to C$ be a cyclic $c$-isogeny, let $\psi := [-1] \circ \varphi \circ \hat{\gamma}\colon C \to E$ with $P_c := \gamma(P_0)$,
$Q_c := \gamma(Q_0)$, so that $\psi(P_c) = -cP$ and $\psi(Q_c) = -cQ$. If we apply the Weil pairing

$$e_{2^a}\colon E[2^a] \times E[2^a] \to \mu_{2^a}(\mathbb{F}_{p^2})$$

to any $R, S \in E[2^a]$ we obtain $e_{2^a}(x\psi(R), x\psi(S)) = e_{2^a}(R, S)^{x^2 c 3^b} = e_{2^a}(R, S)^{-1}$,
which implies that $[x] \circ \psi_{|C[2^a]}\colon C[2^a] \to E[2^a]$ is an anti-isometry, and that

**(2)** $\langle(P_c, x\psi(P_c)), (Q_c, x\psi(Q_c))\rangle = \langle(P_c, -xcP), (Q_c, -xcQ)\rangle = \langle(P_c, P), (Q_c, Q)\rangle$ is a maximal isotropic subgroup of $(C \times E)[2^a]$.

**Definition**

An anti-isometry $\iota\colon C[N] \to E[N]$ is reducible if $(C \times E)/G$ is a product of elliptic curves, where $G := \langle (P_c, \iota(P_c)), (Q_c, \iota(Q_c)) \rangle \subseteq (C \times E)[N]$ for some $\langle P_c, Q_c \rangle = C[N]$.

**Definition**

An isogeny diamond of order $N$ is a triple $(\psi, H_1, H_2)$ with $\psi\colon C \to E$ separable, $H_1, H_2 \subseteq \ker\psi$, $H_1 \cap H_1 = \{0\}$, $H_1 \times H_2 = \ker\psi$, and $\#H_1 + \#H_2 = N$.

**Theorem (Kani97)**

*Let $(\psi, H_1, H_2)$ be an isogeny diamond of order $N$, let $d := \gcd(\#H_1, \#H_2)$, $n := N/d$, and $k_i := \#H_i/d$. Then $\psi = \psi' \circ [d]$ for some $\psi'\colon C \to E$ and there exists a unique reducible anti-isometry $\iota\colon C[N] \to E[N]$ such that*

$$\iota(k_1 R_1 + k_2 R_2) = \psi'(R_2 - R_1)$$

*for all $R_i \in [n]^{-1} H_i$, and every reducible anti-isometry $C[N] \to E[N]$ arises this way.*

## Applying Kani's theorem

Kani's theorem implies that if (1) holds then $x\psi_{|C[2^a]}$ is a reducible anti-isometry and $(C \times E)/\langle (P_c, P), (Q_c, Q) \rangle$ is a product of elliptic curves. Heuristically, the converse almost always holds for large $p$, which allows one to test whether or not (1) holds.

Suppose the prime factors of $c = 2^a - 3^b$ are all congruent to $1 \bmod 4$ (this happens with probability $\approx 1/\sqrt{a}$). Call such a $c$ good. Then we can write $c = u^2 + 4v^2$ and

$$\gamma_{\text{start}} := [u] + [v] \circ 2i \in \text{End}(E_{\text{start}})$$

has degree $c$ and is easy to compute; we can choose $u, v$ so $\gamma_{\text{start}}$ cyclic. We are given a cyclic $3^\beta$-isogeny $\tau \colon E_0 \to E_{\text{start}}$. Let $\tilde\tau$ be the cyclic isogeny with kernel $\gamma_{\text{start}}(\ker \hat\tau)$. Then $\tilde\tau \circ \gamma_{\text{start}} \circ \tau \colon E_0 \to C$ is a $3^{2\beta}c$-isogeny whose kernel contains $E_0[3^\beta]$. Now let

$$\gamma := \frac{\tilde\tau \circ \gamma_{\text{start}} \circ \tau}{[3^\beta]} \colon E_0 \to C.$$

To compute $\gamma(P_0)$ and $\gamma(Q_0)$ we apply $\tilde\tau \circ \gamma_{\text{start}} \tau$ and multiply by $1/3^\beta$ modulo $2^a$.

# Turning our solution to the decision problem into an algorithm

**(1)** Choose $\beta_1 \geq 1$ minimal so there exists $\alpha_1 \geq 0$ with $c_1 := 2^{a-\alpha_1} - 3^{b-\beta_1}$ good.

**(2)** Write $\varphi = \varphi_1 \circ \kappa_1$ where $\kappa_1 \colon E_0 \to E_1$ is a cyclic $3^{\beta_1}$-isogeny
(there are $3^{\beta_1}$ or $4 \cdot 3^{\beta_1 - 1}$ possible $\kappa_1$ and we expect $\beta_1$ to be small, e.g. $\beta_1 = 2$).

**(3)** Test every possible $\kappa_1$ by solving the decision problem on the inputs

  **(ii)** $E_1 := \kappa_1(E_0)$.

  **(iii)** $P_1 := \kappa_1(2^{\alpha_1} P_0)$, $Q_1 := \kappa_1(2^{\alpha_1} Q_0)$ generating $E_1[2^{a-\alpha_1}]$.

  **(iv)** The cyclic $3^{\beta_1}$-isogeny $\kappa_1 \colon E_1 \to E_0$.

  **(v)** The given codomain $E$ of a secret cyclic $3^{b-\beta_1}$-isogeny $\varphi_1 \colon E_1 \to E$.

  **(vi)** The points $2^{\alpha_1} \varphi(P_1)$ and $2^{\alpha_1} \varphi(P_2)$ generating $E[2^{a-\alpha_1}]$.

Step 3 involves an initial gluing step (see Howe-Leprévost-Poonen) followed by a sequence of Richelot isogenies, followed by a single "$\delta = 0$" test, where $\delta$ is the determinant that shows up in the formula for computing Richelot isogenies.

Heuristically exactly one $\kappa_1$ will pass the test and you can stop as soon as you find it.

## Background on Richelot isogenies

For general background on Richelot isogenies and the $(2,2)$-isogeny graph of PPAS's I recommend Benjamin Smith's thesis, and this paper by Smith and Florit.

Let $C: y^2 = f(x)$ be a genus 2 curve defined by a sextic $f \in k[x]$ with Jacobian $J$. The maximal isotropic subgroups of $J[2]$ (kernels of $(2,2)$-isogenies) are in one-to-one correspondence with quadratic splittings of $f$. If $f = g_1 g_2 g_3$ with $g_1, g_2, g_3 \in \bar{k}[x]$ quadratic, the divisors $D_1$ and $D_2$ formed by taking the difference of the points whose $x$-coordinates are roots of $g_1(x)$ and $g_2(x)$ (respectively) generate a maximal isotropic subgroup $G_{g_1 g_2} \subseteq J[2]$, and every maximal isotropic subgroup of $J[2]$ arises this way.

Let $g_i(x) = g_{i2} x^2 + g_{i1} x + g_{i0}$. We may assume $g_{12} = g_{22} = 1$. If

$$\delta := \det \begin{pmatrix} g_{10} & g_{11} & 1 \\ g_{20} & g_{21} & 1 \\ g_{30} & g_{31} & g_{32} \end{pmatrix}$$

is nonzero then $J/G_{g_1 g_2}$ is the Jacobian $J'$ of a genus 2 curve $C'$, and there are explicit formulas for the $(2,2)$-isogeny $J \to J'$ with kernel $G_{g_1 g_2}$ and the curve $C'$. But if $\delta = 0$ then $J/G_{g_1 g_2}$ is a product of elliptic curves; this is the "$\delta = 0$" test.

# Computing all the ternary digits of $sk_B$

Once we know $\kappa_1$ we have the first few digits of $sk_B$. We now choose $\beta_2 > \beta_1$ and proceed as above to compute $\kappa_2$, and then do the same for $\kappa_3, \ldots, \kappa_r$.

In total we need to compute approximately

$$\frac{1}{2}\left(3^{\beta_1} + 3^{\beta_2 - \beta_1} + \cdots + 3^{b - \beta_r}\right)$$

chains of $(2,2)$-isogenies. In the best case this is about $9b/4$, and this is actually close to the typical case, at least heuristically; most steps involve two ternary digits of $sk_B$.

As an optimization, after the first step one can extend $\varphi$ by composing with an extra 3-isogeny, making it easier to obtain good $c_i$ using smaller $\beta_i$. With this optimization we are typically computing just one ternary digit of $sk_B$ in each step.

<p style="text-align:center; color:purple;">Demonstration time!</p>