# An update on lattice cryptanalysis vol. 2

The cost of sieving, and the security margin of Kyber 768

John Schanck Mozilla March 24, 2024

# <u>Concrete cryptanalysis</u> — bit security paradigm



Secure systems take more than 2<sup>128</sup> turns of the crank to break.

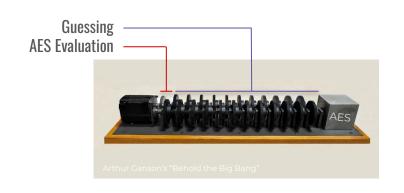
# Resource realism

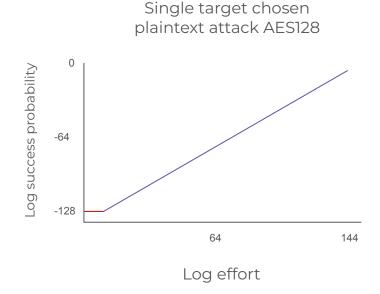
- Real attackers are constrained by:
  - Chip area (mass, system fits on earth),
  - Time (human scale),
  - Power (solar flux, other natural resources),
  - Physical law (locality, finiteness, reliability),
- Real attackers maximize their success probability subject to their constraints.

# max Pr[success | attack, constraints]

-log<sub>2</sub> of this is an operational definition of "security margin"

# <u>Concrete cryptanalysis — resource realist paradigm</u>





Shape of the effort → success probability curve matters. It defines "security margin" for various constraints.

# The dual attack



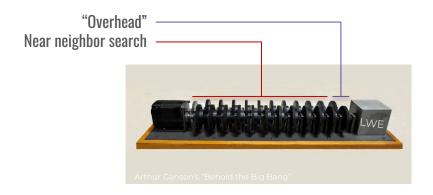
https://github.com/malb/lattice-estimator commit 00ec72c. dual\_hybrid. MATZOV reduction model.

# **Outline**

Lattice attacks on Kyber have:

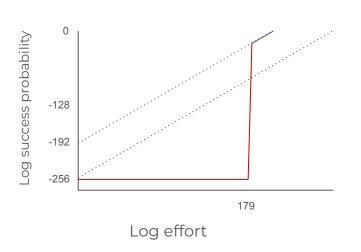
- poor effort → success probability scaling,
- 2. which gets worse when the attacker is memory constrained,
- 3. and even worse when we factor in data movement costs.

## Poor effort → success probability scaling



Small number of iterations Huge cost per-iteration

### Dual attack on Kyber768 using 2-sieve

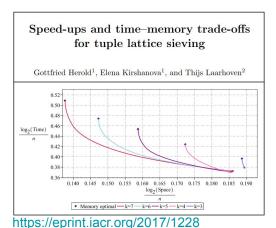


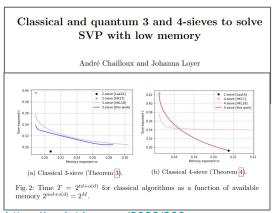
# <u>Outline</u>

Lattice attacks on Kyber have:

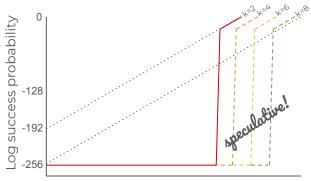
- poor effort → success probability scaling,
- 2. which gets worse when the attacker is memory constrained,
- 3. and even worse when we factor in data movement costs.

### Memory constraints lead to worse effort → success probability scaling





Dual attack on Kyber768 using k-sieve



https://eprint.iacr.org/2023/200

Log effort

For Kyber768 (d=538), I suspect you need a memory exponent  $\sim$ 0.16 (k = 8?) if you are constrained to  $\sim$ 2<sup>100</sup> bits.

# Fermi approximation: Is the attacker memory constrained?

### Facts:

- Industry consumed ~2<sup>43</sup> mm<sup>2</sup> of wafers in 2022.
- 3D NAND density is  $\sim 2^{34}$  bits/mm<sup>2</sup>, or  $2^{43}$  bits/g.  $2^{43}$  mm<sup>2</sup> ·  $2^{34}$  bits / mm<sup>2</sup> =  $2^{77}$  bits.
- The moon has a mass of  $2^{86}$  g.

$$2^{86} g \cdot 2^{43}$$
 bits  $/ g = 2^{129}$  bits.

### Conclusion: Yes.

Need density-production product to scale by 2<sup>50</sup> to store the 2<sup>127</sup> bit database needed for a 2-sieve attack on Kyber768.



### Annual Silicon\* Industry Trends

	2019	2020	2021	2022	2023
Area Shipments (MSI)	11,810	12,407	14,165	14,713	12,602
Revenues (\$Billion)	11.2	11.2	12.6	13.8	12.3

Source: SEMI (www.semi.org), February 2024

\*Data cited in this release include polished silicon wafers, including those used as virgin test wafers, as well as epitaxial silicon wafers, and non-polished silicon wafers shipped by the wafer manufacturers to end users. Shipments are for semiconductor applications only and do not include solar applications.





# <u>Outline</u>

Lattice attacks on Kyber have:

- 1. poor effort → success probability scaling,
- 2. which gets worse when the attacker is memory constrained,
- 3. and even worse when we factor in data movement costs.

# Does memory-access add exponential cost?

Subject of intense discussion for 2-sieves. More work needed for k-sieves



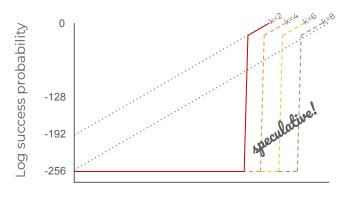
Memory adds no cost to lattice sieving for computers in 3 or more spatial dimensions

Samuel Jaques

Department of Combinatorics and Optimization
University of Waterloo
sejaques@uwaterloo.ca

https://eprint.iacr.org/2024/080

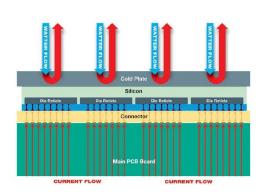
Dual attack on Kyber768 using k-sieve

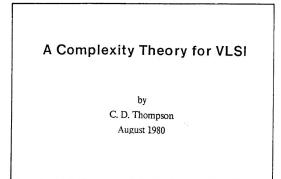


Log effort

Current understanding: all curves move right by (small) exponential factor on 2D mesh architecture.

# Why consider 2D mesh computers?







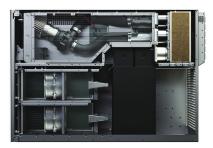


Figure 4: This side view shows the water movement assembly (top), and the air movement infrastructure — fans and a heat exchanger (bottom half).

In particular, the area A and time T taken by any

VLSI chip using any algorithm to perform an N-point Fourier transform must satisfy  $AT^2 \ge cN^2log^2N$ , for some fixed c>0. A more general result for both sorting and Fourier transformation is that  $AT^{2x} = \Omega(N^{1+x}log^{2x}N)$ , for any x in the range  $0 \le x \le I$ . Also, the energy dissipated by a VLSI chip during the solution of either of these problems is at least  $\Omega(N^{3/2}logN)$ . The tightness of these bounds is demonstrated by the existence of nearly optimal circuits for both sorting and Fourier transformation. The circuits based on the shuffle-exchange interconnection pattern are fast but large:  $T = O(log^2N)$  for Fourier transformation,  $T = O(log^3N)$  for sorting; both have area A of at most  $O(N^2/log^{1/2}N)$ . The circuits based on the mesh interconnection pattern are slow but small:  $T = O(N^{1/2}loglogN)$ ,  $A = O(N log^2N)$ .

# Fermi approximation: Cost of memory with mesh routing

### Facts:

- 2022 silicon wafer supply → 2<sup>27.5</sup> WSE-3s
  - $\circ$  2<sup>47.5</sup> cores,
  - o 2<sup>66</sup> bits of memory,
  - 2<sup>85</sup> bits/s mesh bandwidth,
  - Sort 2<sup>66</sup> bits of small data in a few minutes,
  - 4 TW of power.
- Annual global electricity supply ~30000 TWh 30000 TWh / 3600 s = 8.3 TW

### Conclusion:

 Already energy and chip-area constrained for a 2<sup>66</sup> bit mesh sort. Factor 2<sup>61</sup> away from Kyber768 2-sieve size.



### **Annual Silicon\* Industry Trends**

	2019	2020	2021	2022	2023
Area Shipments (MSI)	11,810	12,407	14,165	14,713	12,602
Revenues (\$Billion)	11.2	11.2	12.6	13.8	12.3

Source: SEMI (www.semi.org), February 2024

"Data cited in this release include polished silicon wafers, including those used as virgin test wafers, as well as epitaxial silicon wafers, and non-polished silicon wafers shipped by the wafer manufacturers to end users. Shipments are for semiconductor applications only and do not include solar applications.

### Breakdown of global electricity supply and emissions, 2021-2026

TWh	2021	2022	2023	2026	Growth rate 2021- 2022	Growth rate 2022- 2023
Total Generation	28 426	29 124	29 734	32 694	2.5%	2.1%

Cerebras Wafer-Scale Engine **Fabrication process** 5nm Silicon area 46.225mm Transistors 4 Trillion Al-optimized cores 900,000 Memory (on-chip) **44GB** Memory bandwidth 21PB/s Cerebras WSE-3 Fabric bandwidth 4 Trillion Transistors 214Pb/s 46,225 mm<sup>2</sup> Silicon

# <u>Outline</u>

Lattice attacks on Kyber have:

- 1. poor effort → success probability scaling,
- 2. which gets worse when the attacker is memory constrained,
- 3. and even worse when we factor in data movement costs.

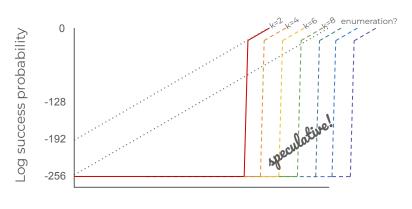
# **Open questions**

- Re-evaluate FFT distinguisher step of the dual attack with memory / interconnect constraints.
- Compute non-asymptotic cost tables for k-sieves.
- Complete "resource realist" analysis of lattice attacks.
  - ... with memory constraints.
  - ... with energy or operation constraints.
- Determine best attack on Kyber768 for constrained adversaries.
  - Seed guessing?
  - Decryption failure attacks?
  - Combinatorial / hybrid attacks?

# <u>Takeaways</u>

- Lattice attacks have poor effort → success scaling.
- 2-sieve memory is unobtainable.
- K-sieving reduces memory but adds exponential cost.
- Interconnect and chip area constraints add further cost, even if only subexponential.
- While there's a significant amount of analysis left to be done, it's not unreasonable to think that Kyber768 is as secure as AES-256.





Log effort