Advanced Lattice Sieving on GPUs, with Tensor Cores

Léo Ducas, Marc Stevens, Wessel van Woerden (CWI).



• Most NIST PQC finalists (5/7) are based on hard lattice problems.

- Most NIST PQC finalists (5/7) are based on hard lattice problems.
- Practical cryptanalysis is important to pick concrete parameters.

- Most NIST PQC finalists (5/7) are based on hard lattice problems.
- Practical cryptanalysis is important to pick concrete parameters.
- Lattice sieving algorithms have the best practical and asymptotic runtime.

- Most NIST PQC finalists (5/7) are based on hard lattice problems.
- Practical cryptanalysis is important to pick concrete parameters.
- Lattice sieving algorithms have the best practical and asymptotic runtime.
- How fit are (different) sieving algorithms for specialized hardware?

- Most NIST PQC finalists (5/7) are based on hard lattice problems.
- Practical cryptanalysis is important to pick concrete parameters.
- Lattice sieving algorithms have the best practical and asymptotic runtime.
- How fit are (different) sieving algorithms for specialized hardware?
- Including more advanced sieving techniques.

- Most NIST PQC finalists (5/7) are based on hard lattice problems.
- Practical cryptanalysis is important to pick concrete parameters.
- Lattice sieving algorithms have the best practical and asymptotic runtime.
- How fit are (different) sieving algorithms for specialized hardware?
- Including more advanced sieving techniques.

Contributions

• First GPU implementation using all state-of-the-art sieving techniques.

- Most NIST PQC finalists (5/7) are based on hard lattice problems.
- Practical cryptanalysis is important to pick concrete parameters.
- Lattice sieving algorithms have the best practical and asymptotic runtime.
- How fit are (different) sieving algorithms for specialized hardware?
- Including more advanced sieving techniques.

Contributions

- First GPU implementation using all state-of-the-art sieving techniques.
- Improves both runtime and energy efficiency by two orders of magnitude.

- Most NIST PQC finalists (5/7) are based on hard lattice problems.
- Practical cryptanalysis is important to pick concrete parameters.
- Lattice sieving algorithms have the best practical and asymptotic runtime.
- How fit are (different) sieving algorithms for specialized hardware?
- Including more advanced sieving techniques.

Contributions

- First GPU implementation using all state-of-the-art sieving techniques.
- Improves both runtime and energy efficiency by two orders of magnitude.
- Significantly improve several lattice problem records.

- Most NIST PQC finalists (5/7) are based on hard lattice problems.
- Practical cryptanalysis is important to pick concrete parameters.
- Lattice sieving algorithms have the best practical and asymptotic runtime.
- How fit are (different) sieving algorithms for specialized hardware?
- Including more advanced sieving techniques.

Contributions

- First GPU implementation using all state-of-the-art sieving techniques.
- Improves both runtime and energy efficiency by two orders of magnitude.
- Significantly improve several lattice problem records.
- First optimized implementation of the asymptotic best known sieve [BDGL].

Lattice



Lattice



2 / 25

Shortest Vector Problem



Bounded Distance Decoding

	•		•		•		•		•		•		•		•		•		•		•		•		•		•	
•		•		•		•		•		•		•		•		•	t	•		•		•		•		•		•
	•		•		•		•		•		•		•		•		•		•		•		•		•		•	
•		•		•		•		•		•		•	0	•		•		•		•		•		•		•		•
	•		•		•		•		•		•		•		•		•		•		•		•		•		•	
•		•		•		•		•		•		•		•		•		•		•		•		•		•		•
	•		•		•		•		•		•		•		•		•		•		•		•		•		•	

Bounded Distance Decoding

	•		•		•		•		•		•		•		•		•		•		•		•		•		•	
•		•		•		•		•		•		•		•		•	+	•		•		•		•		•		•
	•		•		•		•		•		•		•		•		•	5	•		•		•		•		•	
•		•		•		•		•		•		•	(•		•		•		•		•		•		•		•
	•		•		•		•		•		•		•		•		•		•		•		•		•		•	
•		•		•		•		•		•		•		•		•		•		•		•		•		•		•
	•		•		•		•		•		•		•		•		•		•		•		•		•		•	

TU Darmstadt Lattice Challenge

• Gives an indication of the <u>concrete</u> hardness of SVP.

TU Darmstadt Lattice Challenge

- Gives an indication of the <u>concrete</u> hardness of SVP.
- Given: 'Random' *d*-dimensional lattice *L* (Goldstein and Mayer)

TU Darmstadt Lattice Challenge

- Gives an indication of the <u>concrete</u> hardness of SVP.
- Given: 'Random' *d*-dimensional lattice \mathcal{L} (Goldstein and Mayer)
- Goal: Find a $\boldsymbol{v} \in \boldsymbol{\mathcal{L}}$ s.t.

 $\|\mathbf{v}\| \leq 1.05 \cdot \mathsf{GH}(\mathcal{L}) \approx 1.05 \cdot \lambda_1(\mathcal{L}).$

Lattice Sieving



Lattice Sieving



Lattice Sieving



Graphics Processing Unit (GPU)



Graphics Processing Unit (GPU)



Graphics Processing Unit (GPU)





• Very efficient (low precision) matrix multiplication.



- Very efficient (low precision) matrix multiplication.
- 16-bit precision is good enough.



- Very efficient (low precision) matrix multiplication.
- 16-bit precision is good enough.
- Up to 108 16-bit Tflops! [2018 model we used]



- Very efficient (low precision) matrix multiplication.
- 16-bit precision is good enough.
- Up to 108 16-bit Tflops! [2018 model we used]
- Newest model > **300** 16-bit **Tflops**.



- Very efficient (low precision) matrix multiplication.
- 16-bit precision is good enough.
- Up to 108 16-bit Tflops! [2018 model we used]
- Newest model > **300** 16-bit **Tflops**.
- The current best CPU would reach at most ≈ 5 16-bit Tflops.





9 / 25







Hard to adapt algorithms.



<u>Cons</u>

Not versatile.

'External' device.

Memory bottlenecks often limit actual performance.

Hard to adapt algorithms.



<u>Cons</u>

Not versatile.

'External' device.

Memory bottlenecks often limit actual performance.

Hard to adapt algorithms.



G6K, Albrecht et al. 2019

• Open source sieving framework/implementation.
- Open source sieving framework/implementation.
- Advanced sieving: Combines all state-of-the-art results/'tricks'.

- Open source sieving framework/implementation.
- Advanced sieving: Combines all state-of-the-art results/'tricks'.
- Fully parallel (CPU, single machine).

- Open source sieving framework/implementation.
- Advanced sieving: Combines all state-of-the-art results/'tricks'.
- Fully parallel (CPU, single machine).
- \bullet SVP record at dimension 155 in ± 14 days

- Open source sieving framework/implementation.
- Advanced sieving: Combines all state-of-the-art results/'tricks'.
- Fully parallel (CPU, single machine).
- \bullet SVP record at dimension 155 in ± 14 days
 - 4×18 cpu cores.

- Open source sieving framework/implementation.
- Advanced sieving: Combines all state-of-the-art results/'tricks'.
- Fully parallel (CPU, single machine).
- \bullet SVP record at dimension 155 in ± 14 days
 - 4×18 cpu cores.
 - $\approx 256 GiB$ memory.

- Open source sieving framework/implementation.
- Advanced sieving: Combines all state-of-the-art results/'tricks'.
- Fully parallel (CPU, single machine).
- \bullet SVP record at dimension 155 in ± 14 days
 - 4×18 cpu cores.
 - $\approx 256 GiB$ memory.
- Enumeration: dimension 152 using 800.000 core hours!.

Advanced Sieving on GPUs

Sieving Process





Partition the sphere.



Partition the sphere.

Only check all pairs within each bucket.



Partition the sphere.

Only check all pairs within each bucket.

Increases reduction probability per pair.

• No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.
 - \sqrt{N} buckets of size \sqrt{N} .

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.
 - \sqrt{N} buckets of size \sqrt{N} .
 - time $2^{0.349d+o(d)}$.

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.
 - \sqrt{N} buckets of size \sqrt{N} .
 - time $2^{0.349d+o(d)}$.
 - Trick: Replace random directions with random lattice vectors.

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.
 - \sqrt{N} buckets of size \sqrt{N} .
 - time $2^{0.349d+o(d)}$.
 - Trick: Replace random directions with random lattice vectors.
 - Compute inner products with lattice vectors: Tensor cores!

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.
 - \sqrt{N} buckets of size \sqrt{N} .
 - time 2^{0.349d+o(d)}
 - Trick: Replace random directions with random lattice vectors.
 - Compute inner products with lattice vectors: Tensor cores!
- [BDGL] Structured spherical cones (product code).

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.
 - \sqrt{N} buckets of size \sqrt{N} .
 - time $2^{0.349d+o(d)}$.
 - Trick: Replace random directions with random lattice vectors.
 - Compute inner products with lattice vectors: Tensor cores!
- [BDGL] Structured spherical cones (product code).
 - $N^{k/(k+1)}$ buckets of size $N^{1/(k+1)} \leftarrow A$ lot of small buckets!

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.
 - \sqrt{N} buckets of size \sqrt{N} .
 - time $2^{0.349d+o(d)}$.
 - Trick: Replace random directions with random lattice vectors.
 - Compute inner products with lattice vectors: Tensor cores!
- [BDGL] Structured spherical cones (product code).
 - $N^{k/(k+1)}$ buckets of size $N^{1/(k+1)} \leftarrow A$ lot of small buckets!
 - time $2^{0.292d+o(d)}$ for $\mathbf{k} = \theta(\log n)$.

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.
 - \sqrt{N} buckets of size \sqrt{N} .
 - time $2^{0.349d+o(d)}$.
 - Trick: Replace random directions with random lattice vectors.
 - Compute inner products with lattice vectors: Tensor cores!
- [BDGL] Structured spherical cones (product code).
 - $N^{k/(k+1)}$ buckets of size $N^{1/(k+1)} \leftarrow A$ lot of small buckets!
 - time $2^{0.292d+o(d)}$ for $k = \theta(\log n)$.
 - Trick: Implicit directions using permutations and Hadamard transform.

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.
 - \sqrt{N} buckets of size \sqrt{N} .
 - time $2^{0.349d+o(d)}$.
 - Trick: Replace random directions with random lattice vectors.
 - Compute inner products with lattice vectors: Tensor cores!
- [BDGL] Structured spherical cones (product code).
 - $N^{k/(k+1)}$ buckets of size $N^{1/(k+1)} \leftarrow A$ lot of small buckets!
 - time $2^{0.292d+o(d)}$ for $\mathbf{k} = \theta(\log n)$.
 - Trick: Implicit directions using permutations and Hadamard transform.
 - Suitable for CPU (AVX2) and GPU.

- No bucketing: time $N^2 = (4/3)^{d+o(d)} = 2^{0.415d+o(d)}$.
- [BGJ1] Random spherical cones.
 - \sqrt{N} buckets of size \sqrt{N} .
 - time $2^{0.349d+o(d)}$.
 - Trick: Replace random directions with random lattice vectors.
 - Compute inner products with lattice vectors: Tensor cores!
- [BDGL] Structured spherical cones (product code).
 - $N^{k/(k+1)}$ buckets of size $N^{1/(k+1)} \leftarrow A$ lot of small buckets!
 - time $2^{0.292d+o(d)}$ for $\mathbf{k} = \theta(\log n)$.
 - Trick: Implicit directions using permutations and Hadamard transform.
 - Suitable for CPU (AVX2) and GPU.
 - AVX2 CPU implementation merged into G6K, fastest CPU sieve.

Bucketing Quality



• For each pair v, w in a bucket check if $||v \pm w|| < C$.

- For each pair $\boldsymbol{v}, \boldsymbol{w}$ in a bucket check if $\|\boldsymbol{v} \pm \boldsymbol{w}\| < \boldsymbol{C}$.
- $\|\mathbf{v} \pm \mathbf{w}\|^2 = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 \pm 2\langle \mathbf{v}, \mathbf{w} \rangle.$

- For each pair v, w in a bucket check if $||v \pm w|| < C$.
- $\|\mathbf{v} \pm \mathbf{w}\|^2 = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 \pm 2\langle \mathbf{v}, \mathbf{w} \rangle$.
- Need to compute pairwise inner products $\langle v, w \rangle$: Tensor cores!.

- For each pair v, w in a bucket check if $||v \pm w|| < C$.
- $\|\mathbf{v} \pm \mathbf{w}\|^2 = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 \pm 2\langle \mathbf{v}, \mathbf{w} \rangle$.
- Need to compute pairwise inner products $\langle v, w \rangle$: Tensor cores!.
- Sparse output: only return successful pairs.

- For each pair $\boldsymbol{v}, \boldsymbol{w}$ in a bucket check if $\|\boldsymbol{v} \pm \boldsymbol{w}\| < \boldsymbol{C}$.
- $\|\mathbf{v} \pm \mathbf{w}\|^2 = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 \pm 2\langle \mathbf{v}, \mathbf{w} \rangle$.
- Need to compute pairwise inner products $\langle v, w \rangle$: Tensor cores!.
- Sparse output: only return successful pairs.
- FLOP: $O(d \cdot B^2)$, data: $O(d \cdot B)$, ratio improves for larger bucket size B.

Amortizing data throughput



Amortizing data throughput



• Small buckets are memory bound.

Amortizing data throughput



- Small buckets are memory bound.
- Large buckets to reach optimal performance.

16 / 25

BGJ1 vs BDGL



Dimensions for Free [Duc18]

• Sieve in a projected sublattice $\pi_I(\mathcal{L})$ (projected away from first I basis vectors).



Dimensions for Free [Duc18]

- Sieve in a projected sublattice $\pi_I(\mathcal{L})$ (projected away from first I basis vectors).
- Lift the database back to the full lattice [Babai lifting].



Dimensions for Free [Duc18]

- Sieve in a projected sublattice $\pi_I(\mathcal{L})$ (projected away from first I basis vectors).
- Lift the database back to the full lattice [Babai lifting].



• Finds the shortest vector for $I = O\left(\frac{d}{\log(d)}\right)$.
Dimensions for Free [Duc18]

- Sieve in a projected sublattice $\pi_I(\mathcal{L})$ (projected away from first I basis vectors).
- Lift the database back to the full lattice [Babai lifting].



- Finds the shortest vector for $I = O\left(\frac{d}{\log(d)}\right)$.
- Progressive sieving: decrease *I* step-by-step.

Dimensions for Free [Duc18]

- Sieve in a projected sublattice $\pi_I(\mathcal{L})$ (projected away from first I basis vectors).
- Lift the database back to the full lattice [Babai lifting].



- Finds the shortest vector for $I = O\left(\frac{d}{\log(d)}\right)$.
- Progressive sieving: decrease *I* step-by-step.
- On the fly lifting: lift any shortish vector we encounter.

Dimensions for Free [Duc18]

- Sieve in a projected sublattice $\pi_I(\mathcal{L})$ (projected away from first I basis vectors).
- Lift the database back to the full lattice [Babai lifting].



- Finds the shortest vector for $I = O\left(\frac{d}{\log(d)}\right)$.
- Progressive sieving: decrease *I* step-by-step.
- On the fly lifting: lift any shortish vector we encounter.
- Can we efficiently detect if $v_i v_j$ might lift to a short vector [BDD problem]?





19 / 25



19 / 25

• For short dual vectors $z_1,\ldots,z_k\in\mathcal{D}$ we define the **dual hash**

• For short dual vectors $z_1, \ldots, z_k \in \mathcal{D}$ we define the **dual hash**

 $\mathcal{H}(t) := (\langle t, z_i \rangle)_i.$

• If dist (\mathcal{L}, t) is small, then dist $(\mathbb{Z}^k, \mathcal{H}(t))$ is small.

• For short dual vectors $z_1,\ldots,z_k\in\mathcal{D}$ we define the dual hash

- If dist (\mathcal{L}, t) is small, then dist $(\mathbb{Z}^k, \mathcal{H}(t))$ is small.
- For l = 16, k = 48 seems enough for a strong correlation.

• For short dual vectors $z_1,\ldots,z_k\in\mathcal{D}$ we define the dual hash

- If dist (\mathcal{L}, t) is small, then dist $(\mathbb{Z}^k, \mathcal{H}(t))$ is small.
- For l = 16, k = 48 seems enough for a strong correlation.
- $H(t_i t_j) = H(t_i) H(t_j)$, so O(k) operations per pair.

• For short dual vectors $z_1,\ldots,z_k\in\mathcal{D}$ we define the dual hash

- If dist (\mathcal{L}, t) is small, then dist $(\mathbb{Z}^k, \mathcal{H}(t))$ is small.
- For l = 16, k = 48 seems enough for a strong correlation.
- $H(t_i t_j) = H(t_i) H(t_j)$, so O(k) operations per pair.
- Suitable for GPUs.



21 / 25

• G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^{d} \mathbf{x}_i \mathbf{b}_i$:

G6K stores lot's of information per vector y = ∑_{i=1}^d x_ib_i:
[x₁,..., x_d]. (2d)

- G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^{d} \mathbf{x}_i \mathbf{b}_i$:
 - $[x_1,\ldots,x_d]$ (2d)
 - $[\langle \mathbf{y}, \tilde{\mathbf{b}}_1 \rangle / \| \tilde{\mathbf{b}}_1 \|, \dots, \langle \mathbf{y}, \tilde{\mathbf{b}}_d \rangle / \| \tilde{\mathbf{b}}_d \|].$ (4d)

- G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^{d} x_i \mathbf{b}_i$:
 - $[x_1,\ldots,x_d]$ (2d)
 - $[\langle \mathbf{y}, \tilde{\mathbf{b}}_1 \rangle / \| \tilde{\mathbf{b}}_1 \|, \dots, \langle \mathbf{y}, \tilde{\mathbf{b}}_d \rangle / \| \tilde{\mathbf{b}}_d \|].$ (4d)
 - $\|y\|^2$. (8)

- G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^{d} x_i \mathbf{b}_i$:
 - $[x_1,\ldots,x_d]$ (2d)
 - $[\langle \mathbf{y}, \tilde{\mathbf{b}}_1 \rangle / \|\tilde{\mathbf{b}}_1\|, \dots, \langle \mathbf{y}, \tilde{\mathbf{b}}_d \rangle / \|\tilde{\mathbf{b}}_d\|].$ (4d)
 - $\|y\|^2$. (8)
 - Unique Identifier (8)

- G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^{d} x_i \mathbf{b}_i$:
 - $[x_1,\ldots,x_d]$ (2d)
 - $[\langle \mathbf{y}, \tilde{\mathbf{b}}_1 \rangle / \|\tilde{\mathbf{b}}_1\|, \dots, \langle \mathbf{y}, \tilde{\mathbf{b}}_d \rangle / \|\tilde{\mathbf{b}}_d\|].$ (4d)
 - $\|y\|^2$. (8)
 - Unique Identifier (8)
 - Lift target t. (41)

- G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^{d} x_i \mathbf{b}_i$:
 - $[x_1,\ldots,x_d]$ (2d)
 - $[\langle \mathbf{y}, \tilde{\mathbf{b}}_1 \rangle / \|\tilde{\mathbf{b}}_1\|, \dots, \langle \mathbf{y}, \tilde{\mathbf{b}}_d \rangle / \|\tilde{\mathbf{b}}_d\|].$ (4d)
 - $\|y\|^2$. (8)
 - Unique Identifier (8)
 - Lift target t. (41)
 - Dual Hash (41)

- G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^{d} x_i \mathbf{b}_i$:
 - $[x_1,\ldots,x_d]$ (2d)
 - $[\langle \mathbf{y}, \tilde{\mathbf{b}}_1 \rangle / \|\tilde{\mathbf{b}}_1\|, \dots, \langle \mathbf{y}, \tilde{\mathbf{b}}_d \rangle / \|\tilde{\mathbf{b}}_d\|].$ (4d)
 - $\|y\|^2$. (8)
 - Unique Identifier (8)
 - Lift target t. (41)
 - Dual Hash (41)
 - Popcount (*d*/**4**)

- G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^{d} x_i \mathbf{b}_i$:
 - $[x_1,\ldots,x_d]$ (2d)
 - $[\langle \mathbf{y}, \tilde{\mathbf{b}}_1 \rangle / \|\tilde{\mathbf{b}}_1\|, \dots, \langle \mathbf{y}, \tilde{\mathbf{b}}_d \rangle / \|\tilde{\mathbf{b}}_d\|].$ (4d)
 - $\|y\|^2$. (8)
 - Unique Identifier (8)
 - Lift target *t*. (41)
 - Dual Hash (41)
 - Popcount (*d*/**4**)
- Remove all except x, $||y||^2$ and unique identifier.

- G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^{d} x_i \mathbf{b}_i$:
 - $[x_1,\ldots,x_d]$. (2d)
 - $[\langle \mathbf{y}, \tilde{\mathbf{b}}_1 \rangle / \|\tilde{\mathbf{b}}_1\|, \dots, \langle \mathbf{y}, \tilde{\mathbf{b}}_d \rangle / \|\tilde{\mathbf{b}}_d\|].$ (4d)
 - $\|y\|^2$. (8)
 - Unique Identifier (8)
 - Lift target *t*. (41)
 - Dual Hash (41)
 - Popcount (*d*/4)
- Remove all except x, $||y||^2$ and unique identifier.
- Reduces memory by $\pm 60\%$.

- G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^{d} x_i \mathbf{b}_i$:
 - $[x_1,\ldots,x_d]$. (2d)
 - $[\langle \mathbf{y}, \tilde{\mathbf{b}}_1 \rangle / \|\tilde{\mathbf{b}}_1\|, \dots, \langle \mathbf{y}, \tilde{\mathbf{b}}_d \rangle / \|\tilde{\mathbf{b}}_d\|].$ (4d)
 - $\|y\|^2$. (8)
 - Unique Identifier (8)
 - Lift target t. (41)
 - Dual Hash (41)
 - Popcount (*d*/4)
- Remove all except x, $||y||^2$ and unique identifier.
- Reduces memory by $\pm 60\%$.
- Compute everything on the GPU, overhead of $O(B \cdot d^2)$ for a bucket size **B**.

New SVP records



23 / 25

• Lattice sieving algorithms can efficiently be implemented on GPUs.

- Lattice sieving algorithms can efficiently be implemented on GPUs.
- Memory bottlenecks disappear when buckets are large enough.

- Lattice sieving algorithms can efficiently be implemented on GPUs.
- Memory bottlenecks disappear when buckets are large enough.
- Extra benefit of saving memory with negligible overhead.

- Lattice sieving algorithms can efficiently be implemented on GPUs.
- Memory bottlenecks disappear when buckets are large enough.
- Extra benefit of saving memory with negligible overhead.
- BDGL beats BGJ1 in practice on CPUs, but the cross-over for GPUs lies much higher.

Bibliography

- [BGJ15] A. Becker, N. Gama, A. Joux. Speeding-up lattice sieving without increasing the memory, using sub-quadratic nearest neighbor search.
- [BDGL16] A. Becker, L. Ducas, N. Gama, T. Laarhoven. New directions in nearest neighbor searching with applications to lattice sieving.
- [Duc18] L. Ducas, Shortest vector from lattice sieving: A few dimensions for free.
- [LM18] T. Laarhoven, A. Mariano, Progressive lattice sieving.
- [G6K] M.R. Albrecht, L. Ducas, G. Herold, E. Kirshanova, E.W. Postlethwaite, and M. Stevens, 2019. The general sieve kernel and new records in lattice reduction.