

Advanced Lattice Sieving on GPUs, with Tensor Cores

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



Centrum Wiskunde & Informatica

Overview

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

Overview

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

Overview

- 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.

Overview

- 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?

Overview

- 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.

Overview

- 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.

Overview

- 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.

Overview

- 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.

Overview

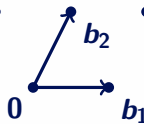
- Most NIST PQC finalists (**5/7**) are based on hard lattice problems.
- Practical cryptanalysis is important to pick concrete parameters.
- Lattice sieving algorithms have the best practical and asymptotic runtime.
- How fit are (different) sieving algorithms for specialized hardware?
- Including more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

$$\mathcal{L} := B\mathbb{Z}^d$$



Lattice

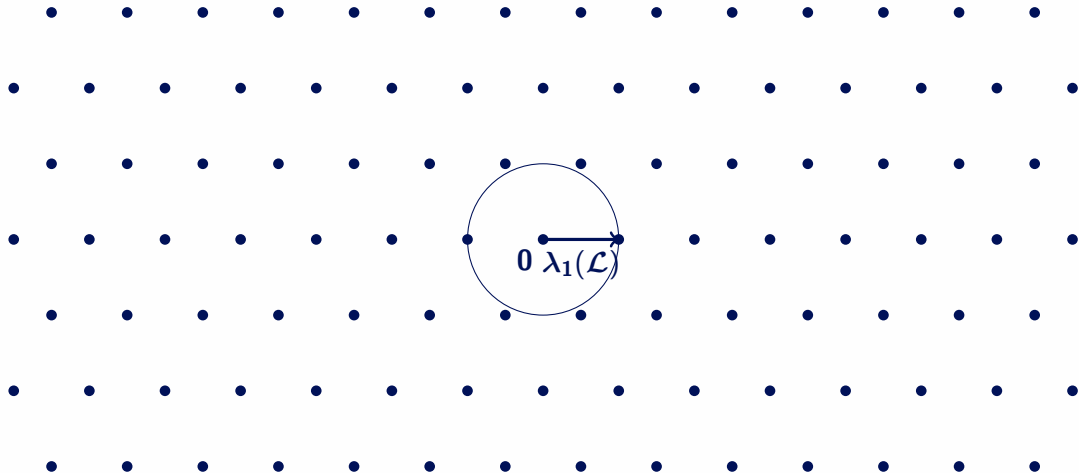
$$\mathcal{L} := B\mathbb{Z}^d$$

0

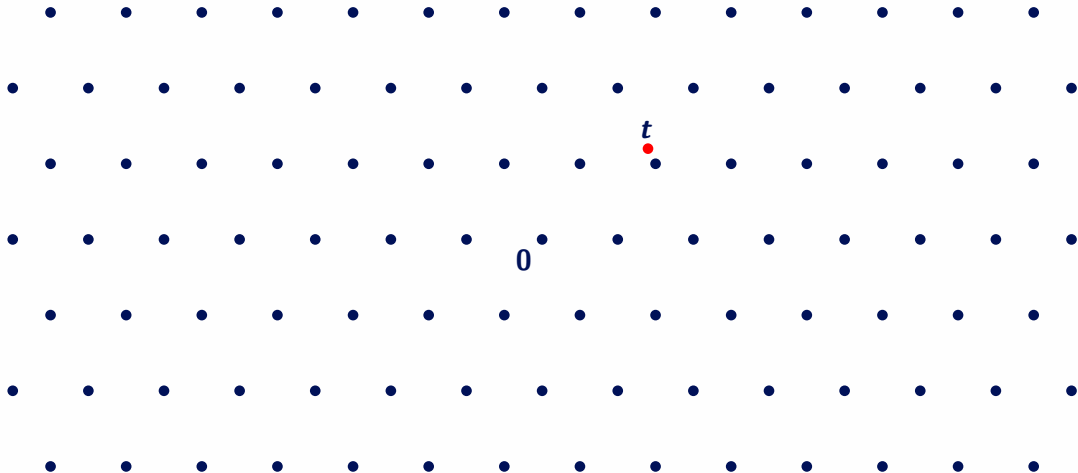
b_2

b_1

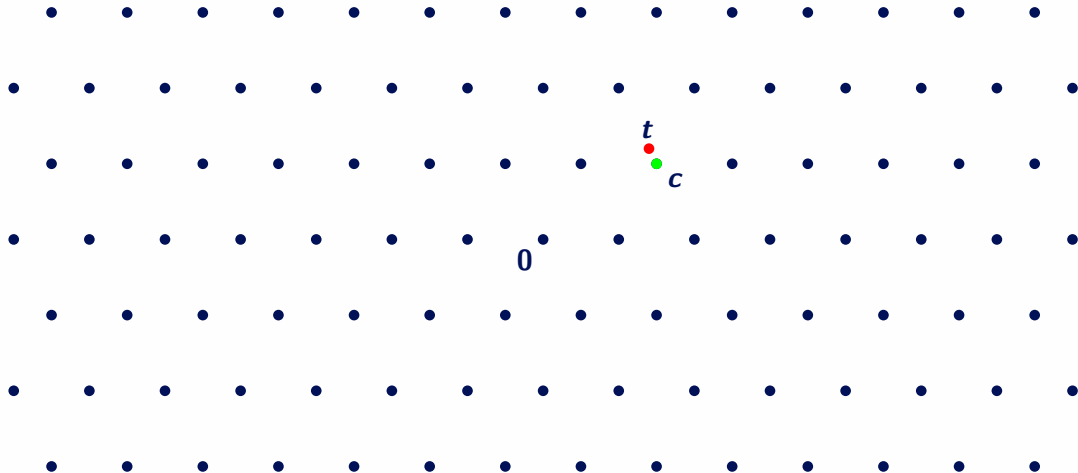
Shortest Vector Problem



Bounded Distance Decoding



Bounded Distance Decoding



TU Darmstadt Lattice Challenge

- Gives an indication of the concrete hardness of SVP.

TU Darmstadt Lattice Challenge

- Gives an indication of the concrete hardness of SVP.
- Given: 'Random' d -dimensional lattice \mathcal{L} (Goldstein and Mayer)

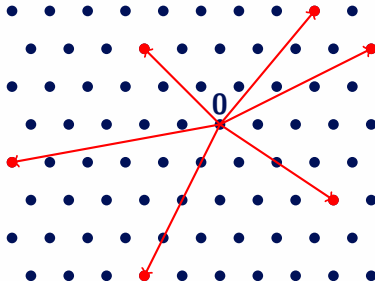
TU Darmstadt Lattice Challenge

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

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

Lattice Sieving

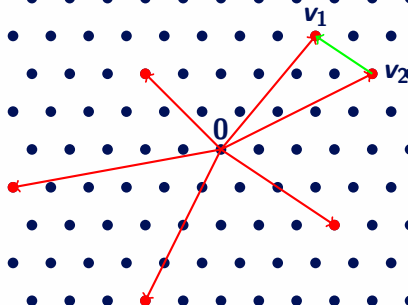
Start with a big list
of $N = (4/3)^{d/2+o(d)}$
lattice vectors.



Lattice Sieving

Start with a big list of $N = (4/3)^{d/2+o(d)}$ lattice vectors.

Find pairs of close vectors to create smaller lattice vectors

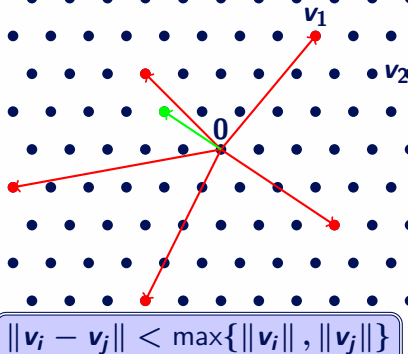


$$\|v_i - v_j\| < \max\{\|v_i\|, \|v_j\|\}$$

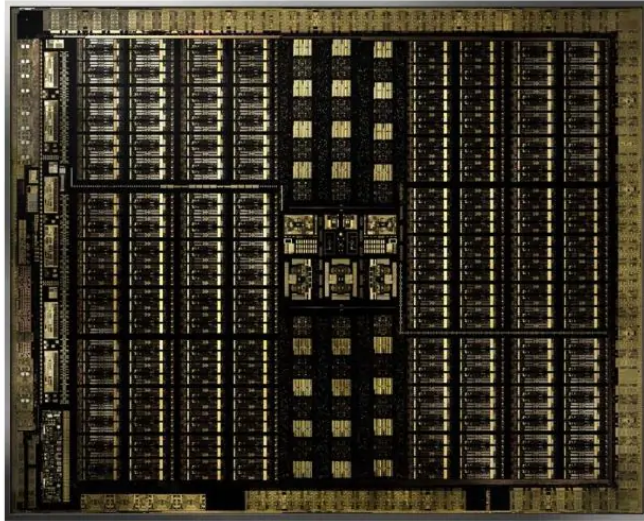
Lattice Sieving

Start with a big list of $N = (4/3)^{d/2+o(d)}$ lattice vectors.

Find pairs of close vectors to create smaller lattice vectors



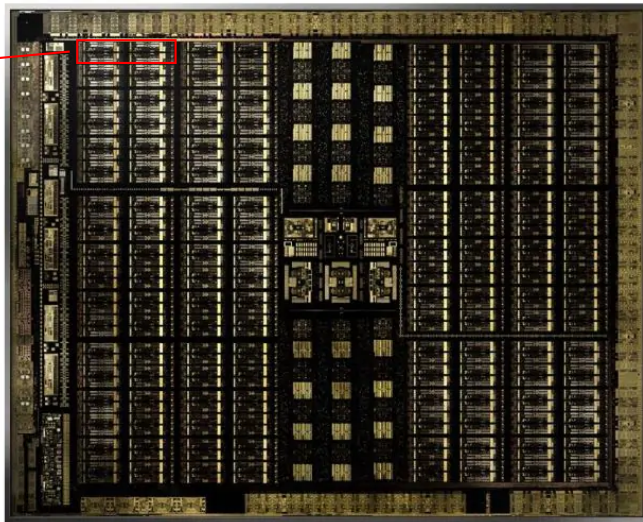
Graphics Processing Unit (GPU)



Graphics Processing Unit (GPU)

64 FP32 cores
64 INT32 cores
8 Tensor cores.

Thousands of
cores.

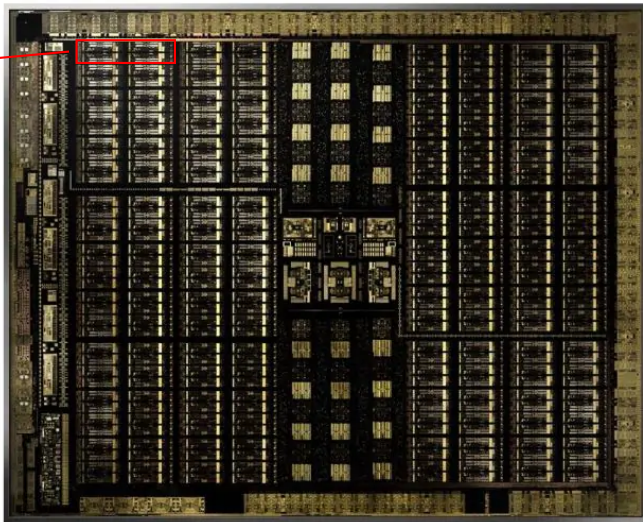


Graphics Processing Unit (GPU)

64 FP32 cores
64 INT32 cores
8 Tensor cores.

Thousands of
cores.

Per **32** cores:
Single Instruction
Multiple Data



Tensor Cores

$$\mathbf{D} = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32 FP16 FP16 or FP32

- Very efficient (low precision) matrix multiplication.

Tensor Cores

$$\mathbf{D} = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32 FP16 FP16 or FP32

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

Tensor Cores

$$\mathbf{D} = \begin{pmatrix} \begin{matrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{matrix} & \begin{matrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{matrix} & \begin{matrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{matrix} \end{pmatrix}$$

FP16 or FP32 FP16 FP16 or FP32

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

Tensor Cores

$$\mathbf{D} = \underbrace{\begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix}}_{\text{FP16 or FP32}} \underbrace{\begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix}}_{\text{FP16}} + \underbrace{\begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}}_{\text{FP16 or FP32}}$$

- 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**.

Tensor Cores

$$\mathbf{D} = \begin{pmatrix} \begin{matrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{matrix} & \begin{matrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{matrix} & \begin{matrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{matrix} \end{pmatrix}$$

FP16 or FP32 FP16 FP16 or FP32

- 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 \approx **5** 16-bit **Tflops**.

GPU: pros and cons

Cons

Not versatile.

Pros

GPU: pros and cons

Cons

Not versatile.

'External' device.

Pros

GPU: pros and cons

Cons

Not versatile.

'External' device.

Memory bottlenecks often
limit actual performance.

Pros

GPU: pros and cons

Cons

Not versatile.

'External' device.

Memory bottlenecks often
limit actual performance.

Hard to adapt algorithms.

Pros

GPU: pros and cons

Cons

Not versatile.

'External' device.

Memory bottlenecks often
limit actual performance.

Hard to adapt algorithms.

Pros



GPU: pros and cons

Cons


Not versatile.

'External' device.

Memory bottlenecks often
limit actual performance.

Hard to adapt algorithms.

Pros



Incredible
performance.



Energy
efficient.

- 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).
- SVP record at dimension **155** in ± 14 days

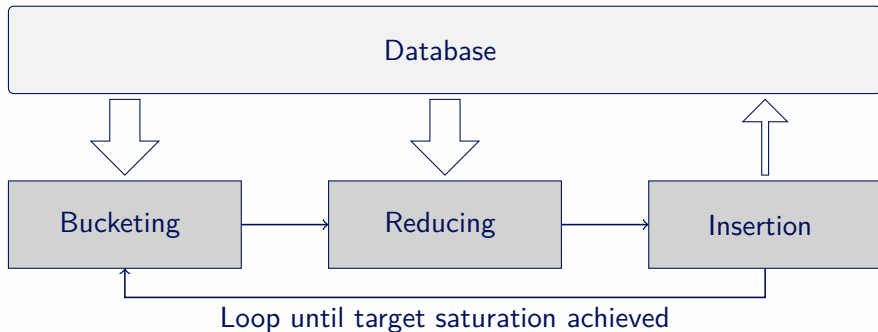
- Open source sieving framework/implementation.
- Advanced sieving: Combines all state-of-the-art results/‘tricks’.
- Fully parallel (CPU, single machine).
- SVP record at dimension **155** in ± 14 days
 - 4×18 cpu cores.

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

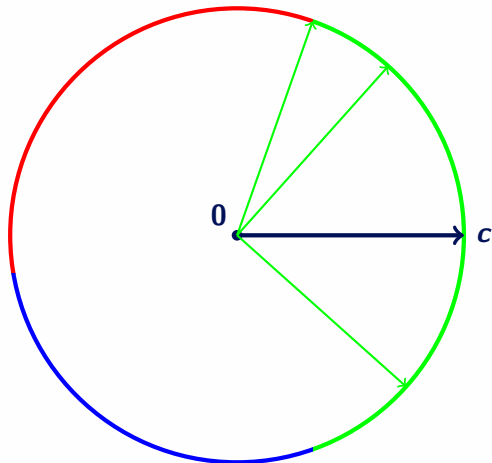
- Open source sieving framework/implementation.
- Advanced sieving: Combines all state-of-the-art results/‘tricks’.
- Fully parallel (CPU, single machine).
- SVP record at dimension **155** in ± 14 days
 - 4×18 cpu cores.
 - ≈ 256 **GiB** memory.
- Enumeration: dimension **152** using **800.000** core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

Advanced Sieving on GPUs

Sieving Process

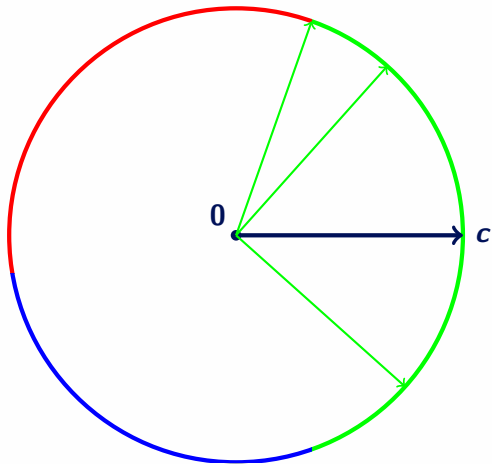


Bucketing



Partition the sphere.

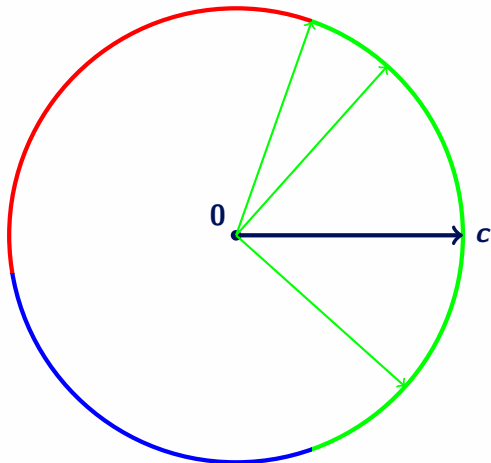
Bucketing



Partition the sphere.

Only check all pairs within each bucket.

Bucketing



Partition the sphere.

Only check all pairs within each bucket.

Increases reduction probability per pair.

Bucketing

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

Bucketing

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

Bucketing

- 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} .

Bucketing

- 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 $2^{0.349d+o(d)}$.

Bucketing

- 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.

Bucketing

- 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!**

Bucketing

- 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).

Bucketing

- 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 $N^{1/(k+1)}$ ← A lot of **small buckets!**

Bucketing

- 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 $N^{1/(k+1)}$ ← A lot of **small buckets!**
 - time $2^{0.292d+o(d)}$ for $k = \theta(\log n)$.

Bucketing

- 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 $N^{1/(k+1)}$ ← A lot of **small buckets!**
 - time $2^{0.292d+o(d)}$ for $k = \theta(\log n)$.
 - Trick: Implicit directions using permutations and Hadamard transform.

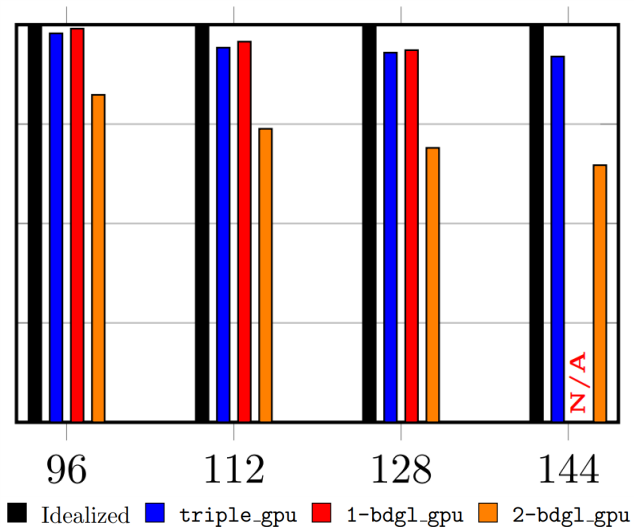
Bucketing

- 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 $N^{1/(k+1)}$ ← A lot of **small buckets!**
 - time $2^{0.292d+o(d)}$ for $k = \theta(\log n)$.
 - Trick: Implicit directions using permutations and Hadamard transform.
 - Suitable for CPU (AVX2) and GPU.

Bucketing

- 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 $N^{1/(k+1)}$ ← A lot of **small buckets!**
 - time $2^{0.292d+o(d)}$ for $k = \theta(\log n)$.
 - 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



Reducing

- For each pair \mathbf{v}, \mathbf{w} in a bucket check if $\|\mathbf{v} \pm \mathbf{w}\| < C$.

Reducing

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

Reducing

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

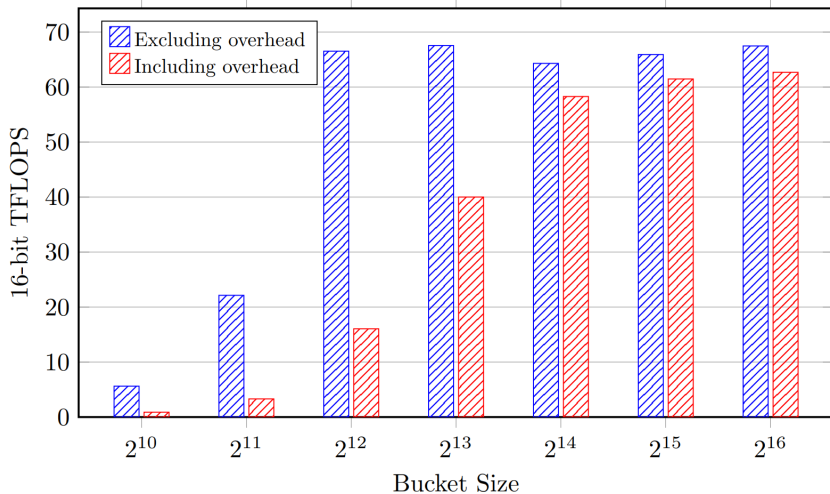
Reducing

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

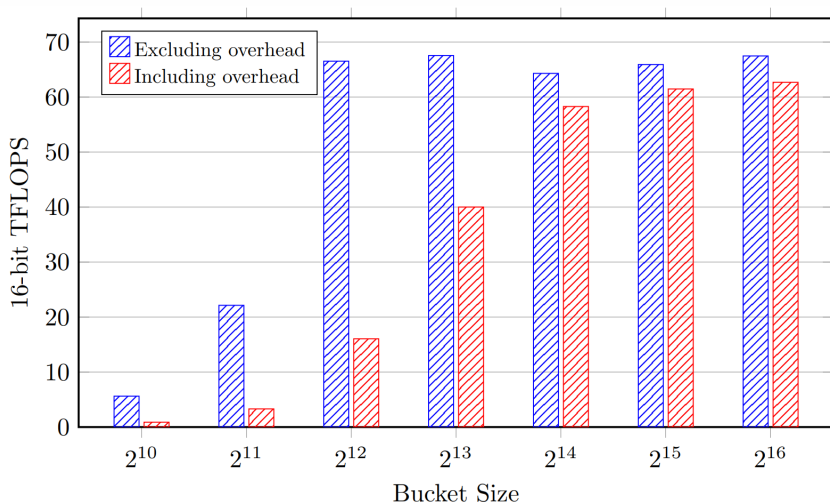
Reducing

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

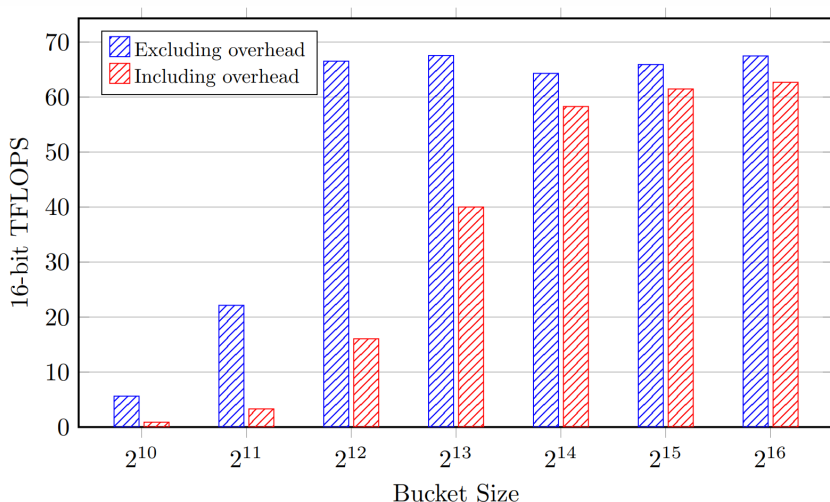


Amortizing data throughput



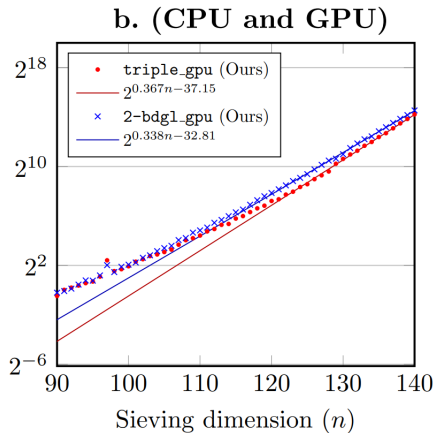
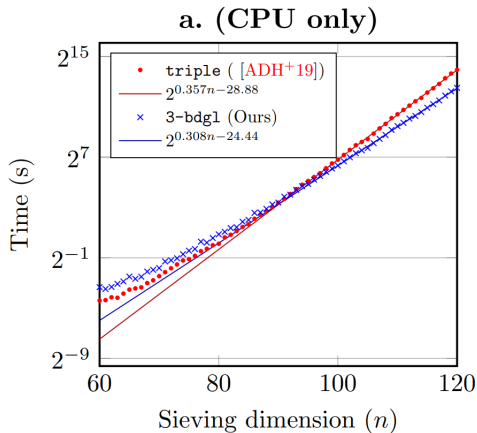
- Small buckets are memory bound.

Amortizing data throughput



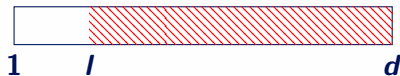
- Small buckets are memory bound.
- Large buckets to reach optimal performance.

BGJ1 vs BDGL



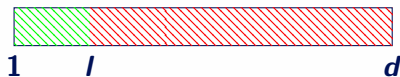
Dimensions for Free [Duc18]

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



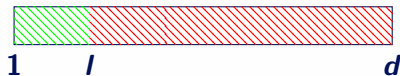
Dimensions for Free [Duc18]

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



Dimensions for Free [Duc18]

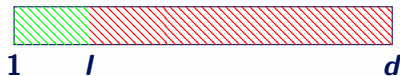
- Sieve in a projected sublattice $\pi_I(\mathcal{L})$ (projected away from first I basis vectors).
- Lift the database back to the full lattice [Babai lifting].



- Finds the shortest vector for $I = O\left(\frac{d}{\log(d)}\right)$.

Dimensions for Free [Duc18]

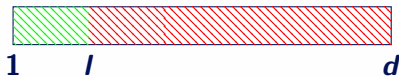
- Sieve in a projected sublattice $\pi_I(\mathcal{L})$ (projected away from first I basis vectors).
- Lift the database back to the full lattice [Babai lifting].



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

Dimensions for Free [Duc18]

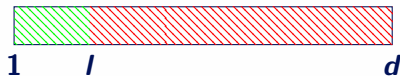
- Sieve in a projected sublattice $\pi_I(\mathcal{L})$ (projected away from first I basis vectors).
- Lift the database back to the full lattice [Babai lifting].



- Finds the shortest vector for $I = \mathcal{O}\left(\frac{d}{\log(d)}\right)$.
- 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 $\mathbf{v}_i - \mathbf{v}_j$ might lift to a short vector [BDD problem]?

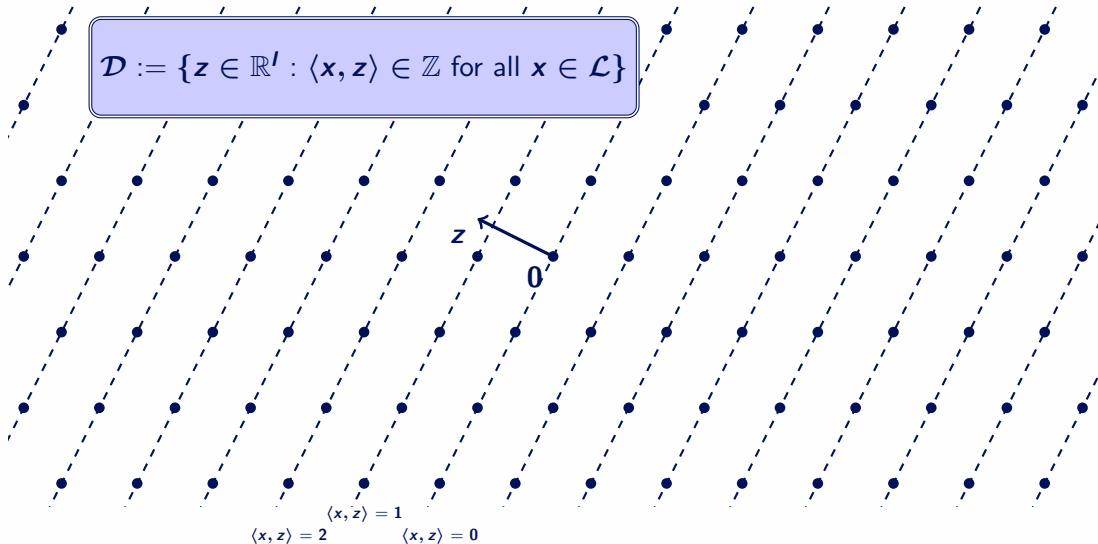
Dual Hash

$$\mathcal{D} := \{z \in \mathbb{R}^I : \langle x, z \rangle \in \mathbb{Z} \text{ for all } x \in \mathcal{L}\}$$

0

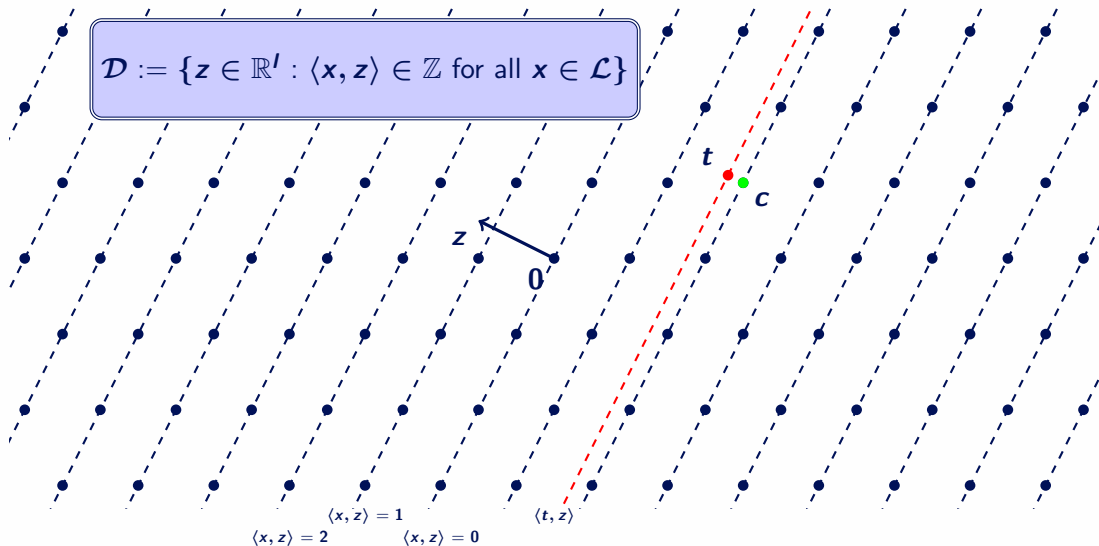
Dual Hash

$$\mathcal{D} := \{z \in \mathbb{R}^I : \langle x, z \rangle \in \mathbb{Z} \text{ for all } x \in \mathcal{L}\}$$



Dual Hash

$$\mathcal{D} := \{z \in \mathbb{R}^I : \langle x, z \rangle \in \mathbb{Z} \text{ for all } x \in \mathcal{L}\}$$



Dual Hash

- For short dual vectors $\mathbf{z}_1, \dots, \mathbf{z}_k \in \mathcal{D}$ we define the **dual hash**

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

Dual Hash

- For short dual vectors $\mathbf{z}_1, \dots, \mathbf{z}_k \in \mathcal{D}$ we define the **dual hash**

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

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

Dual Hash

- For short dual vectors $\mathbf{z}_1, \dots, \mathbf{z}_k \in \mathcal{D}$ we define the **dual hash**

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

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

Dual Hash

- For short dual vectors $\mathbf{z}_1, \dots, \mathbf{z}_k \in \mathcal{D}$ we define the **dual hash**

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

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

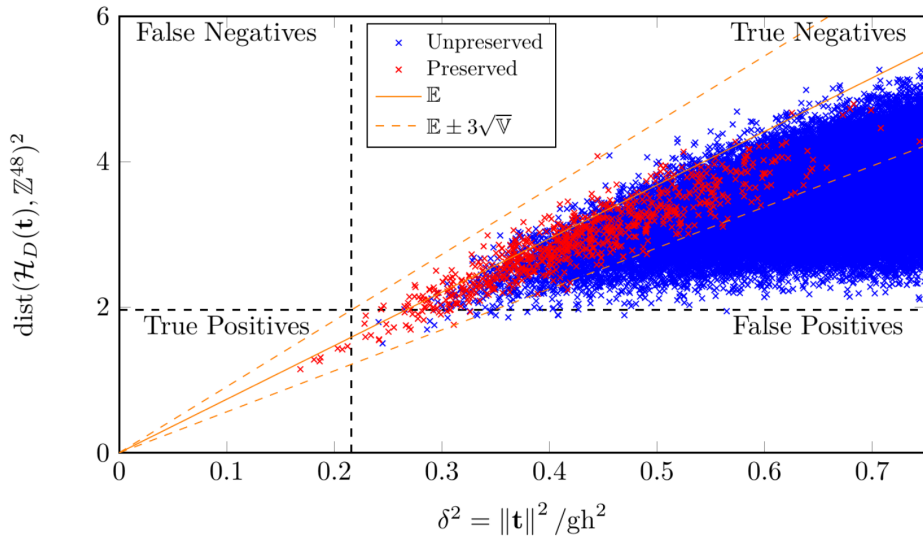
Dual Hash

- For short dual vectors $\mathbf{z}_1, \dots, \mathbf{z}_k \in \mathcal{D}$ we define the **dual hash**

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

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

Dual Hash



Saving Memory

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

Saving Memory

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

Saving Memory

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

Saving Memory

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

Saving Memory

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

Saving Memory

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

Saving Memory

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

Saving Memory

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

Saving Memory

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

Saving Memory

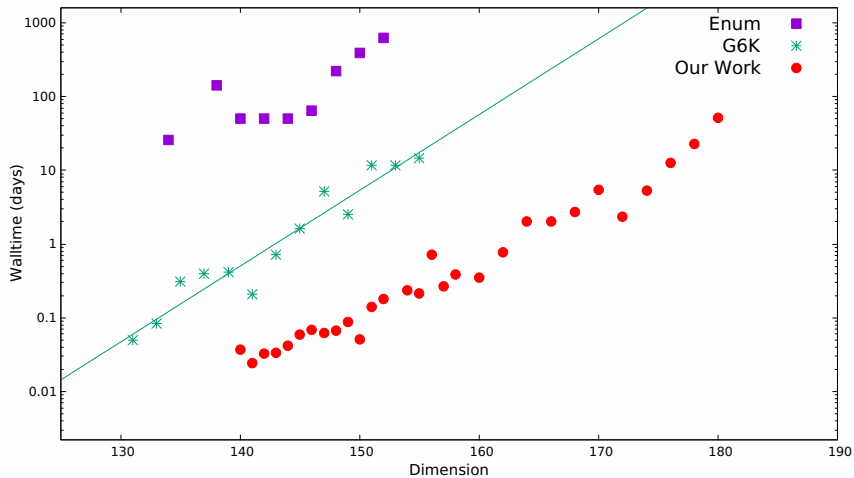
- G6K stores lot's of information per vector $\mathbf{y} = \sum_{i=1}^d \mathbf{x}_i \mathbf{b}_i$:
 - $[\mathbf{x}_1, \dots, \mathbf{x}_d]$. ($2d$)
 - $[\langle \mathbf{y}, \tilde{\mathbf{b}}_1 \rangle / \|\tilde{\mathbf{b}}_1\|, \dots, \langle \mathbf{y}, \tilde{\mathbf{b}}_d \rangle / \|\tilde{\mathbf{b}}_d\|]$.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

Saving Memory

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

Dimension 180!



maximum RAM size of 1.5TB reached for **180**.

Conclusion

- Lattice sieving algorithms can efficiently be implemented on GPUs.

Conclusion

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

Conclusion

- 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.

Conclusion

- 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.