

Shortest Lattice Vector Enumeration on Graphics Cards

Jens Hermans¹ Michael Schneider² Frédéric Vercauteren¹
Johannes Buchmann² Bart Preneel¹

¹K.U.Leuven

²TU Darmstadt

SHARCS - 10 September 2009

- 1 Introduction
- 2 Lattices: crash course
- 3 GPUs
- 4 The Algorithm
- 5 Results
- 6 The Future

Why GPU?



(Source: MSI)

CUDA framework

Warning: sales talk

Your own personal supercomputer for $< \text{€}500$.

Nvidia CUDA Framework:

- Run 'general' programs on GPU
- More complex operations, data types, branching...
- Recent GPU required
- Theory: 1TFlop (practice: 200 GFlop)

Crypto on GPU

Current applications:

- Ciphers:
 - RSA ¹, ECC ², AES ³
- Cryptanalysis:
 - Factoring ⁴
 - Brute force

Focus: high throughput, not latency

¹Moss, Page, Smart / Szerwinski, Guneyasu / Fleissner

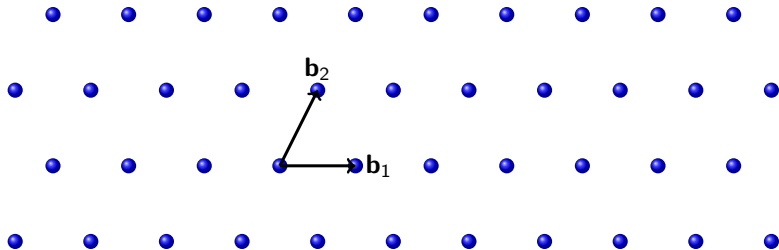
²Szerwinski, Guneyasu

³Manavski / Harrison, Waldron

⁴Bernstein, Chen, Cheng, Lange, Yang

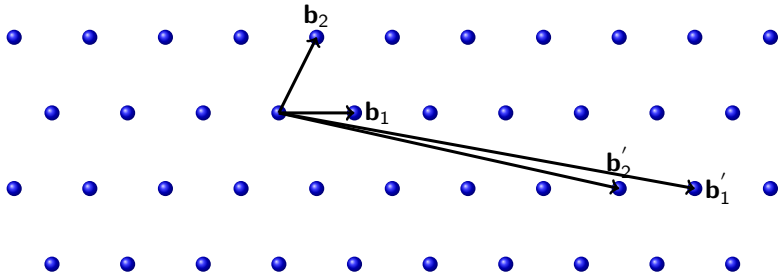
- 1 Introduction
- 2 Lattices: crash course
- 3 GPU
- 4 The Algorithm
- 5 Results
- 6 The Future

Lattices



- Basis matrix $\mathbf{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_n\}$ with $\mathbf{b}_i \in \mathbb{R}^d$
- Lattice: $L(\mathbf{B}) = \{\sum_{i=1}^n x_i \mathbf{b}_i, x_i \in \mathbb{Z}\}$

Shortest Vector Problem (SVP)



- Basis not unique
- Idea: 'good' basis \mathbf{B} and 'bad' basis \mathbf{B}'
- Finding $\lambda_1(L)$ is hard with \mathbf{B}'

Algorithms for SVP

Shortest vector problem

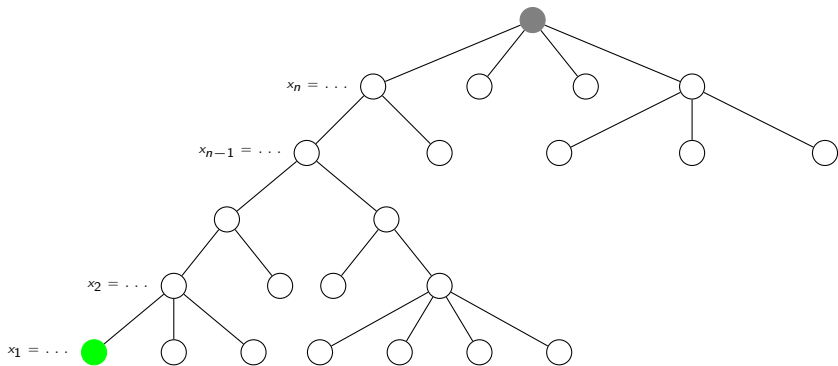
$$\text{Compute } \min_{\mathbf{x} \in \mathbb{Z}^n} \|\mathbf{B}\mathbf{x}\|_2$$

SVP algorithms:

- LLL (+variants): approximate solution, polynomial
- BKZ
- ...
- Enum: exact solution, exponential

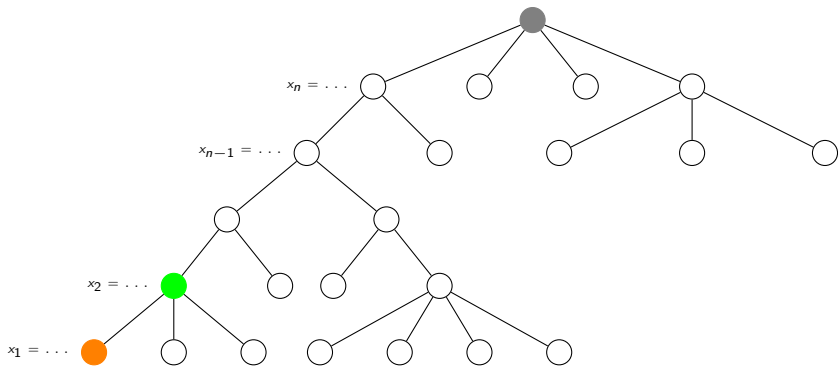
⇒ This talk: focus on enum.

Enumeration



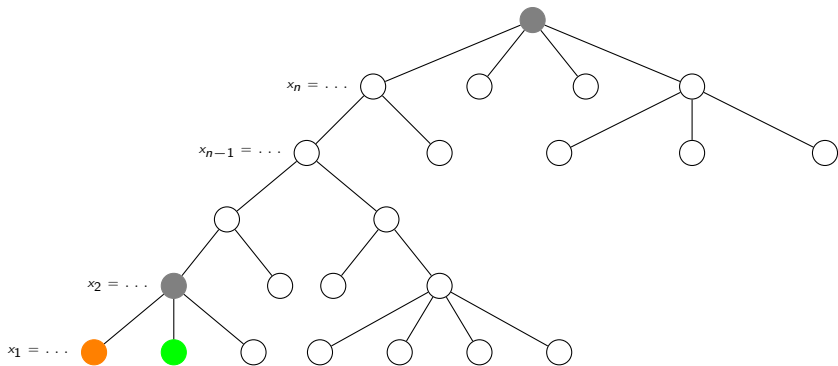
Optimum $A = \|\mathbf{B}\mathbf{x}\|_2^2$ and $\mathbf{x} = [1, 0, \dots, 0]$

Enumeration



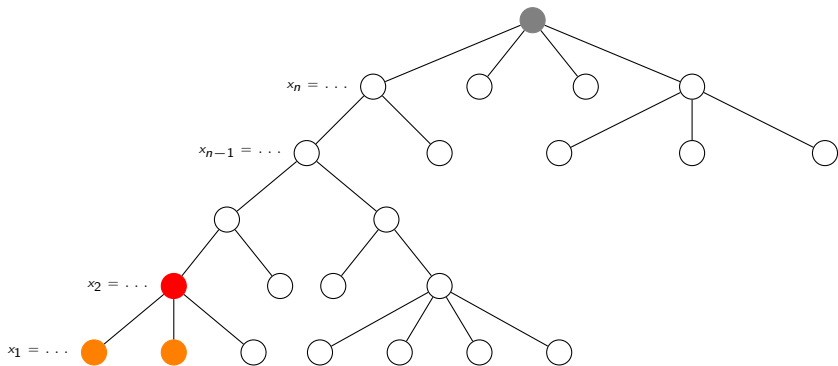
Intermediate norm l_2 s.t. $l_i \geq l_{i+1}$ (with $l_1 = \|\mathbf{B}\mathbf{x}\|_2^2$)

Enumeration



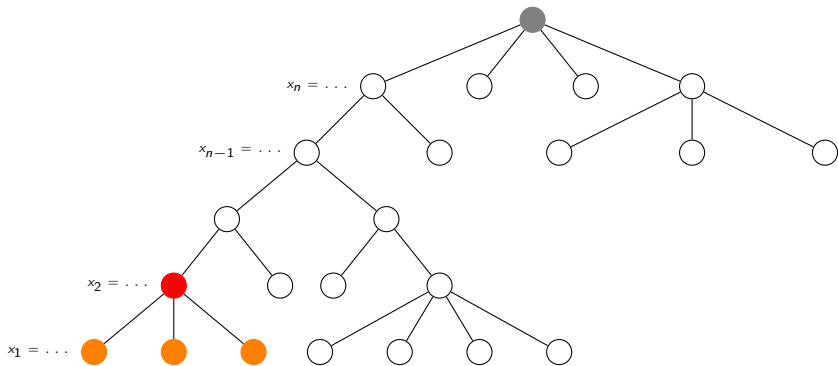
New optimum $A = \|\mathbf{B}\mathbf{x}\|_2^2$

Enumeration

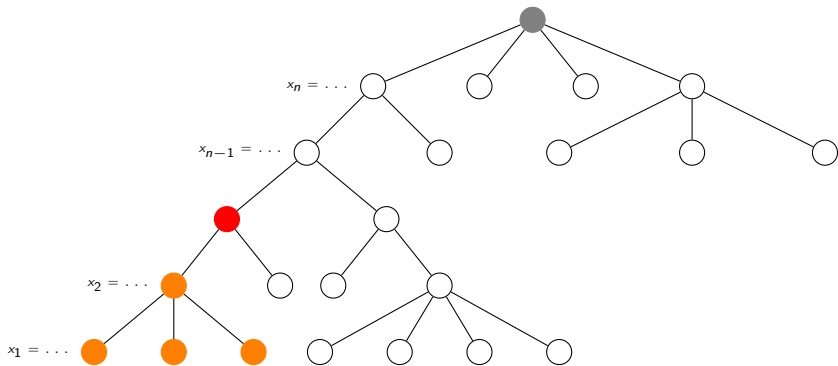


Cut off branch if $l_i > A$.

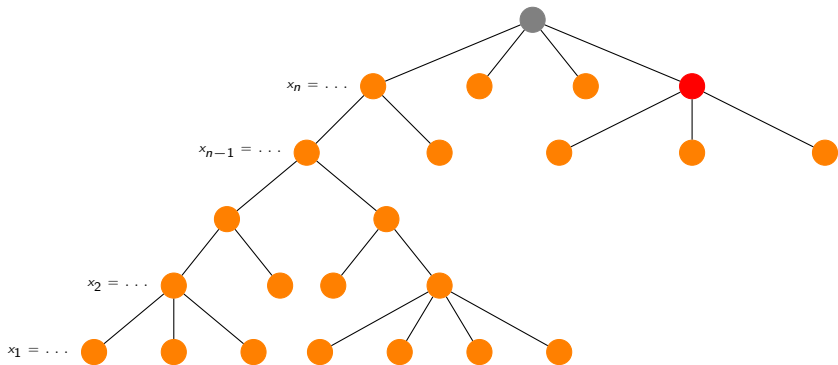
Enumeration



Enumeration



Enumeration



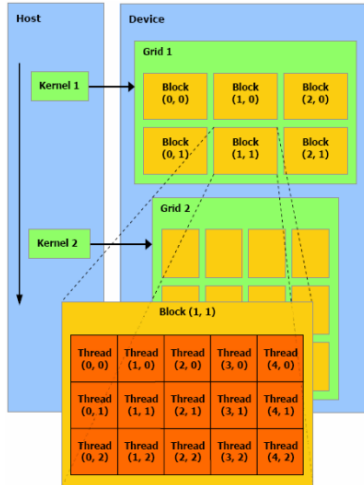
- 1 Introduction
- 2 Lattices: crash course
- 3 GPUs**
- 4 The Algorithm
- 5 Results
- 6 The Future

Processor

Nvidia GTX280:

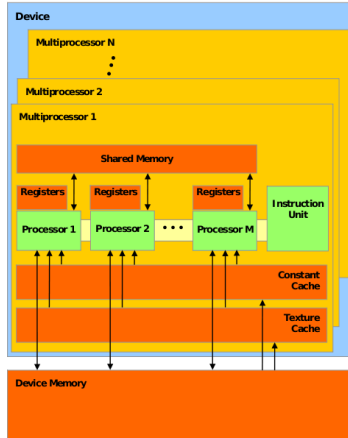
- 240 cores, scalar processors
- 30 multiprocessors (8 cores each)
- 1.3 GHz
- 1GB Global Memory
- 32 & 64-bit integers, FP

Programming model



(Source: CUDA programming guide)

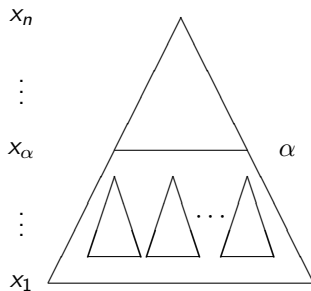
Memory types



(Source: CUDA programming guide)

- 1 Introduction
- 2 Lattices: crash course
- 3 GPUs
- 4 The Algorithm**
- 5 Results
- 6 The Future

Algorithm Flow



Basic idea

Input: \mathbf{B}, A, α, n

- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2
- 3

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$

Basic idea

Input: \mathbf{B} , A , α , n

- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2 CPU: generate $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 3

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$

Basic idea

Input: \mathbf{B} , A , α , n

- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2 CPU: generate $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 3 GPU thread: run a sub-enum on \mathbf{x}_i , if new optimum \rightarrow store in \mathbf{x}

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$

Basic idea

Input: \mathbf{B} , A , α , n

- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2 CPU: generate $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 3 GPU thread: run a sub-enum on \mathbf{x}_i , if new optimum \rightarrow store in \mathbf{x}

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$

\implies horrible performance

Early termination...

Input: \mathbf{B} , A , α , n

- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2 CPU: generate $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$

8

9

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$

Early termination...

Input: \mathbf{B}, A, α, n

- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2 CPU: generate $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 3 GPU thread:
- 4 **while** *there are \mathbf{x}_i left.* **do**
- 5 | Start enum for a certain $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 6 | Stop enum after \mathbf{S} steps, store the state $\{l_i, \bar{\mathbf{x}}_i, s_i = \mathbf{S}\}$
- 7 **end**
- 8
- 9

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$

Early termination...

Input: \mathbf{B}, A, α, n

- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2 CPU: generate $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 3 GPU thread:
- 4 **while** *there are \mathbf{x}_i left.* **do**
- 5 | Start enum for a certain $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 6 | Stop enum after \mathbf{S} steps, store the state $\{l_i, \bar{\mathbf{x}}_i, s_i = \mathbf{S}\}$
- 7 **end**
- 8 CPU: Get enum state $\bar{\mathbf{x}}_i = [\bar{x}_1, \dots, \bar{x}_{\alpha-1}, x_\alpha, \dots, x_n]$
- 9

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$

Early termination...

Input: \mathbf{B} , A , α , n

- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2 CPU: generate $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 3 GPU thread:
- 4 **while** *there are \mathbf{x}_i left.* **do**
- 5 | Start enum for a certain $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 6 | Stop enum after \mathbf{S} steps, store the state $\{l_i, \bar{\mathbf{x}}_i, s_i = \mathbf{S}\}$
- 7 **end**
- 8 CPU: Get enum state $\bar{\mathbf{x}}_i = [\bar{x}_1, \dots, \bar{x}_{\alpha-1}, x_\alpha, \dots, x_n]$
- 9 CPU: Continue enum if \mathbf{S} was reached

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$

\implies solves length difference problem, still not so good

Iterating

Input: \mathbf{B} , A , α , n

- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2 **while** *true* **do**
- 3 CPU: generate some $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 9 CPU: Get enum state $\bar{\mathbf{x}}_i = [\bar{x}_1, \dots, \bar{x}_{\alpha-1}, x_\alpha, \dots, x_n]$
- 10 **end**

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$

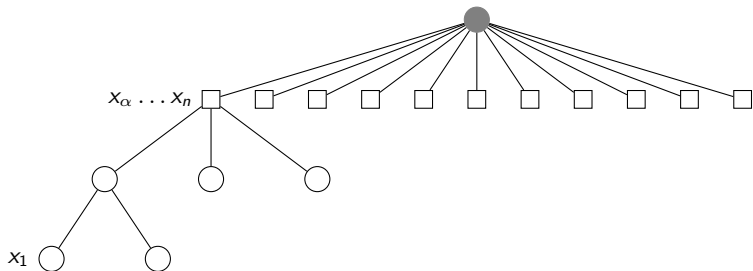
Iterating

Input: \mathbf{B} , A , α , n

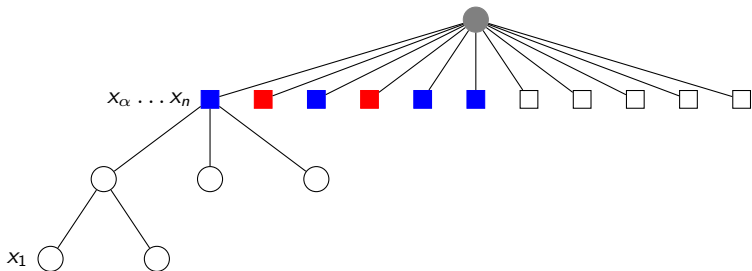
- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2 **while** *true* **do**
- 3 CPU: generate some $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$
- 4 GPU thread:
- 5 **while** *there are \mathbf{x}_i left.* **do**
- 6 Start enum for a certain \mathbf{x}_i or continue enum for $\bar{\mathbf{x}}_i$
- 7 Stop enum after \mathbf{S} steps, store the state $\{i, \bar{\mathbf{x}}_i, s_i = \mathbf{S}\}$
- 8 **end**
- 9 CPU: Get enum state $\bar{\mathbf{x}}_i = [\bar{x}_1, \dots, \bar{x}_{\alpha-1}, x_\alpha, \dots, x_n]$
- 10 **end**

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$

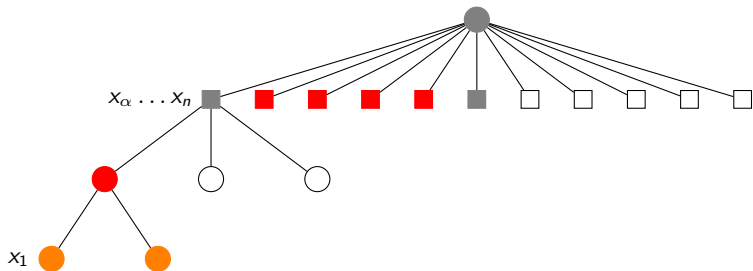
GPU Enumeration



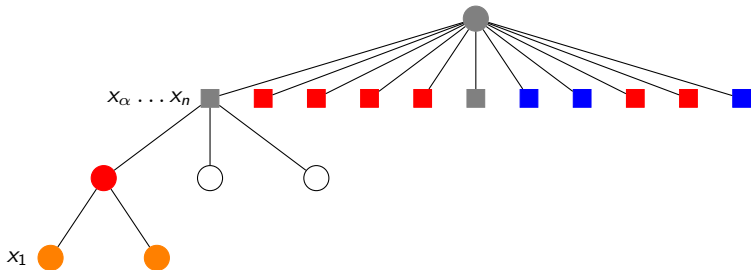
GPU Enumeration



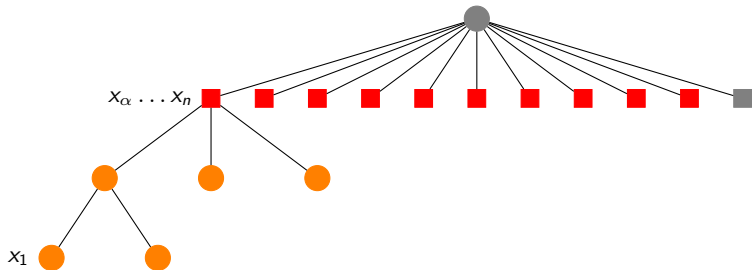
GPU Enumeration



GPU Enumeration



GPU Enumeration



Implementation details

Some facts & figures:

- Dimension 50, 100000 starting vectors \rightarrow upload & download \sim 20 MB of data to GPU
- CPU top enum: very fast (low dimension)
- GPU runs for $>$ 10 seconds per iteration, iteration overhead is limited
- Share new optimal values among GPU threads

- 1 Introduction
- 2 Lattices: crash course
- 3 GPUs
- 4 The Algorithm
- 5 Results**
- 6 The Future

Throughput

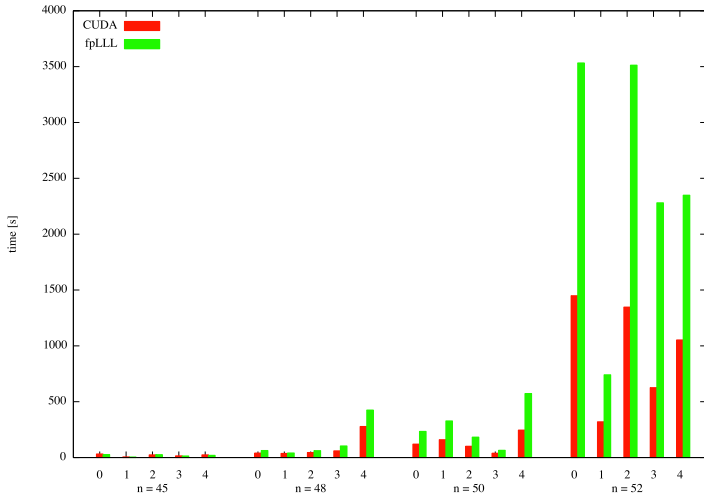
Throughput:

- CPU: around $25 \cdot 10^6$ steps/s
- GPU: up to $100 \cdot 10^6$ steps/s

Throughput on GPU depends on:

- Lattice dimension n
- Length of sub-enumerations
- Number of parallel threads, uploaded points...

Results



Results

n	45	48	50	52	54
fpLLL	18.3s	139s	277s	2483s	6960s
CUDA	20.2s	92s	133s	959s	2599s
	110%	66%	48%	39%	37%

Table: Average time needed for enumeration of lattices in each dimension n .

Ideas for the Future

Future:

- Generalize ideas (not specific for gpu's... clusters?)
- Use full power of CPU (now: idle during gpu-time)
- Gaussian heuristic

The end...

Questions?

Algorithm

Algorithm 1: High-level GPU ENUM Algorithm

Input: $\mathbf{b}_i, A, \alpha, n$

- 1 Compute the Gram-Schmidt decomposition of \mathbf{b}_i
- 2 **while** *true* **do**
- 3 $S = \{(\mathbf{x}_i, \Delta \mathbf{x}_i, \Delta^2 \mathbf{x}_i, l_i = \alpha, s_i = 0)\}_i \leftarrow$ Top enum: generate at most
 NUMSTARTPOINTS $- \#T$ vectors
- 4 $R = \{(\bar{\mathbf{x}}_i, \Delta \mathbf{x}_i, \Delta^2 \mathbf{x}_i, l_i, s_i)\}_i \leftarrow$ GPU enumeration, starting from $S \cup T$
- 5 $T \leftarrow \{R_i : s_i \geq \mathbf{S}\}$
- 6 **if** $\#T < \text{CPUTHRESHOLD}$ **then**
- 7 Enumerate the starting points in T on the CPU.
- 8 Stop
- 9 **end**
- 10 **end**

Output: (x_1, \dots, x_n) with $\|\sum_{i=1}^n x_i \mathbf{b}_i\| = \lambda_1(L)$
