

# Parallel Shortest Lattice Vector Enumeration on Graphics Cards

Jens Hermans<sup>1</sup>

Johannes Buchmann<sup>2</sup>

Bart Preneel<sup>1</sup>

Michael Schneider<sup>2</sup>

Frederik Vercauteren<sup>1</sup>

<sup>1</sup>K.U.Leuven

<sup>2</sup> TU Darmstadt

[mischnei@cdc.informatik.tu-darmstadt.de](mailto:mischnei@cdc.informatik.tu-darmstadt.de)



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT





Introduction

Lattices: crash course

GPUs

GPU-Enum Algorithm

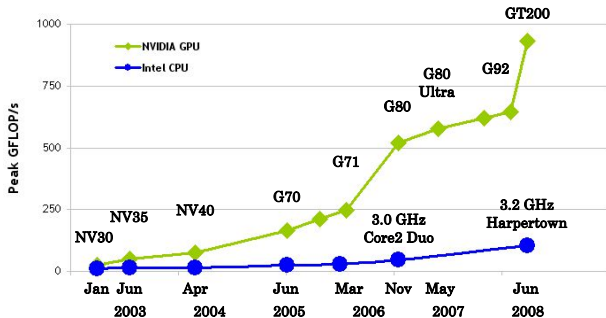
The Future

# Why GPU?



(Source: MSI)

# Why GPU?



GT200 = GeForce GTX 280	G71 = GeForce 7900 GTX	NV35 = GeForce FX 5950 Ultra
G92 = GeForce 9800 GTX	G70 = GeForce 7800 GTX	NV30 = GeForce FX 5800
G80 = GeForce 8800 GTX	NV40 = GeForce 6800 Ultra	

(Source: NVIDIA)

# Why GPU?

Absolute transistor count:

AMD Athlon 64 X2	CPU	154 mio.
Intel Core 2 Duo	CPU	291 mio.
Intel Core 2 Quad	CPU	582 mio.
NVIDIA G8800 GTX	GPU	680 mio.
Intel Core i7-980X Extreme	CPU	1.2 mrd.
ATI RV870	GPU	2.2 mrd.
NVIDIA Fermi	GPU	3.0 mrd.

## Warning: Sales Talk

Your own personal supercomputer for  $< €500$ .

Nvidia CUDA Framework:

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

## Current Crypto Applications

- Ciphers:
  - RSA<sup>a</sup>, ECC<sup>b</sup>, AES<sup>c</sup>
- Cryptanalysis:
  - Factoring<sup>d</sup>
  - Brute force



---

<sup>a</sup>Moss, Page, Smart / Szerwinski, Guneyusu / Fleissner

<sup>b</sup>Szerwinski, Guneyusu

<sup>c</sup>Manavski / Harrison, Waldron

<sup>d</sup>Bernstein, Chen, Cheng, Lange, Yang



Introduction

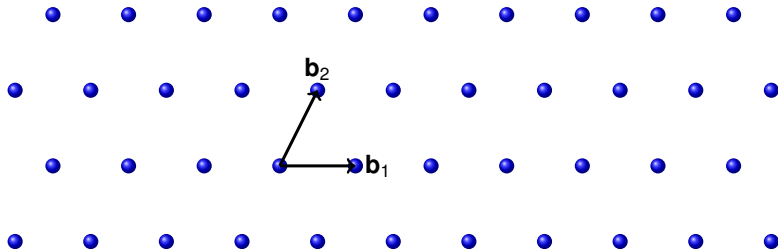
**Lattices: crash course**

GPUs

GPU-Enum Algorithm

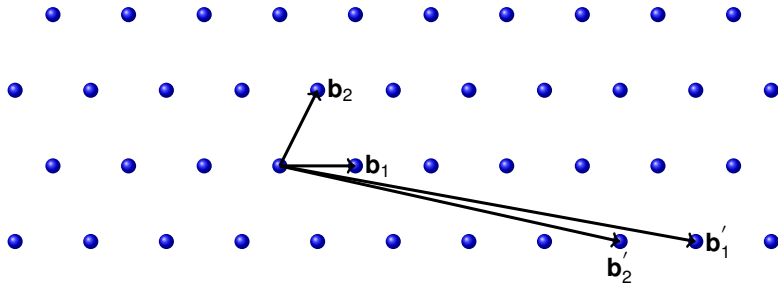
The Future





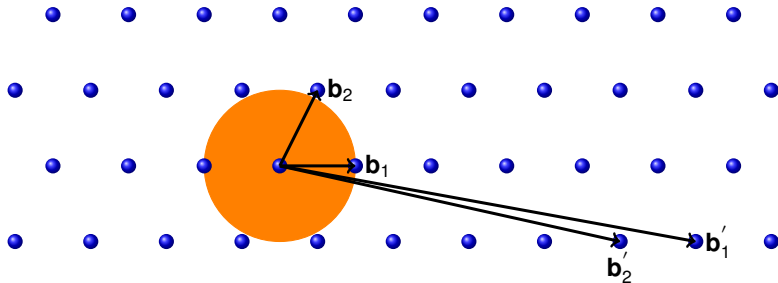
- Basis matrix  $\mathbf{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_n\}$  with  $\mathbf{b}_i \in \mathbb{R}^d$
- Lattice:  $\mathcal{L}(\mathbf{B}) = \{\sum_{i=1}^n x_i \mathbf{b}_i, x_i \in \mathbb{Z}\}$

# Shortest Vector Problem (SVP)



- Basis not unique
- Idea: 'good' basis  $\mathbf{B}$  and 'bad' basis  $\mathbf{B}'$
- Finding a shortest vector is hard with  $\mathbf{B}'$

# Shortest Vector Problem (SVP)



- Basis not unique
- Idea: 'good' basis  $\mathbf{B}$  and 'bad' basis  $\mathbf{B}'$
- Finding a shortest vector is hard with  $\mathbf{B}'$

## Shortest Vector Problem

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

SVP Algorithms:

**LLL** approximate solution, polynomial time

**BKZ** better approximation

...

**Enum** exact solution, exponential time

**Sieving** exact solution, exponential time, probabilistic

⇒ This talk: focus on **Enum**



Introduction

Lattices: crash course

**GPUs**

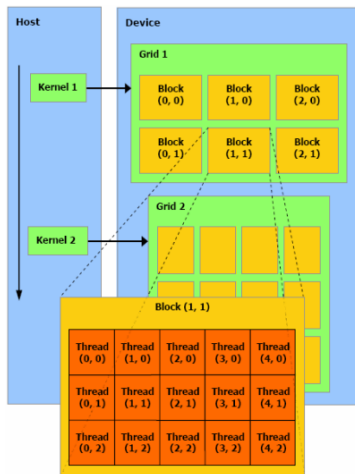
GPU-Enum Algorithm

The Future

## Nvidia GTX280:

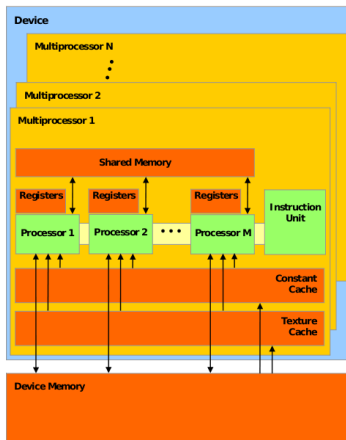
- 240 cores, scalar processors
- 30 multiprocessors (8 cores each)
- 1.3 GHz
- 1GB Global Memory
- 32 & 64-bit integers, single + double precision FP

# Programming model



(Source: CUDA programming guide)

# Memory types



(Source: CUDA programming guide)



## Code Example



```
CUDA_SAFE_CALL(cudaMalloc((void **)&d_A, sizeof(norminttype)));  
CUDA_SAFE_CALL(cudaMemcpy(d_A,h_A,sizeof(norminttype),cudaMemcpyHostToDevice));  
CUDA_SAFE_CALL(cudaFree(d_A));  
gpuenum <<<16,32>>>(paramlist);
```

## Code Example

```
CUDA_SAFE_CALL(cudaMalloc((void **)&d_A, sizeof(norminttype)));
CUDA_SAFE_CALL(cudaMemcpy(d_A,h_A,sizeof(norminttype),cudaMemcpyHostToDevice));
CUDA_SAFE_CALL(cudaFree(d_A));
gpuenum <<<16,32>>>(paramlist);

__global__ void gpuenum(paramlist)
{
    for (int k = 0; k < alpha + 1; k++) {
        fptype cintsumtmp = 0.0;
        for (int j = alpha; j <= beta; j++) {
            cintsumtmp += (x[j + sharedOffset] * MU(j, k));
        }
    }
    if (startlevel >= alpha && l > localA)
        continue;
}
```



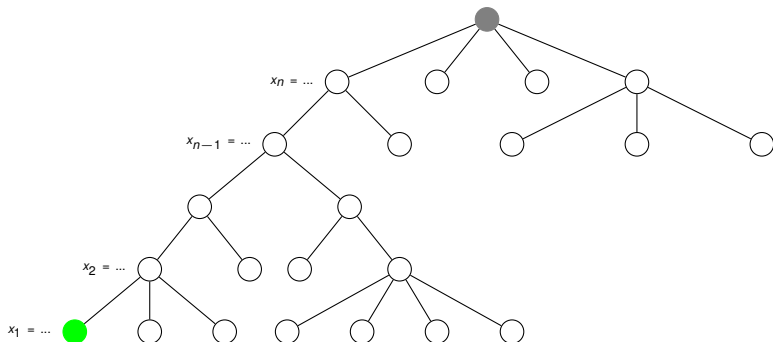
Introduction

Lattices: crash course

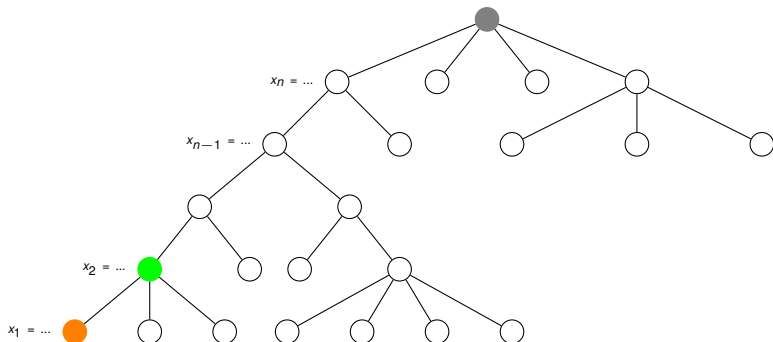
GPUs

**GPU-Enum Algorithm**

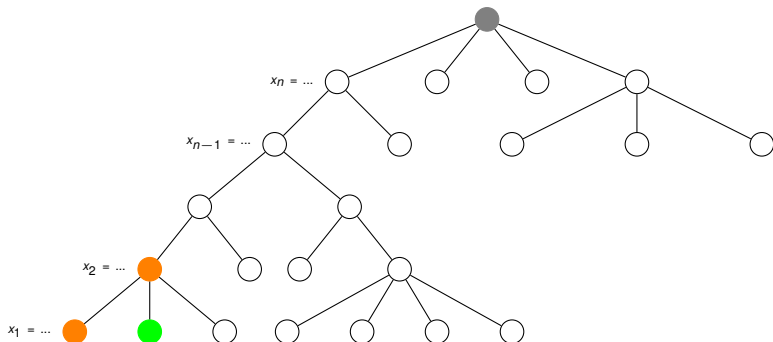
The Future



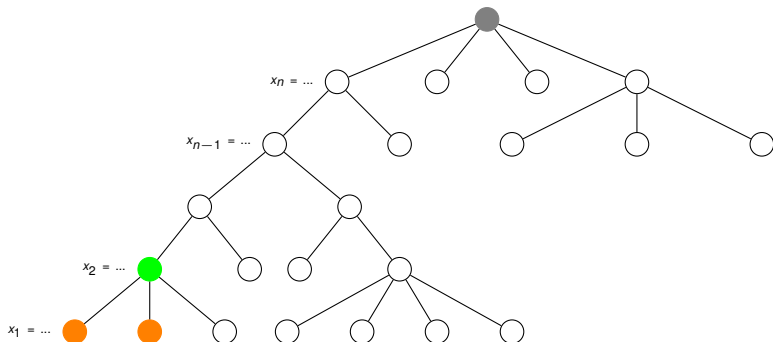
$$\mathbf{x} = [1, 0, \dots, 0] \rightarrow \text{optimum } A = \|\mathbf{B}\mathbf{x}\|_2^2$$



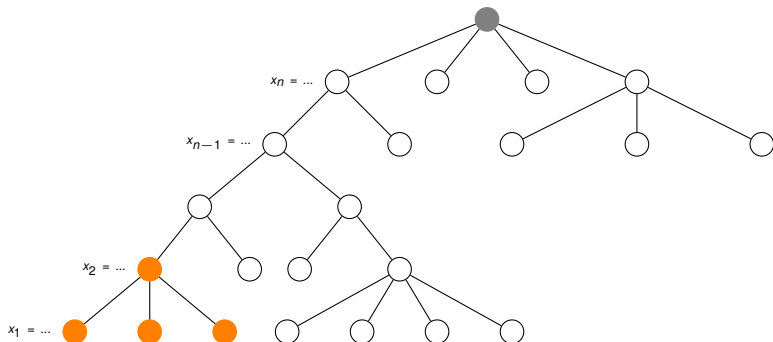
Intermediate norm  $l_2$  of vector  $[\star, x_2, x_3, \dots, x_n]$   
with  $l_i \geq l_{i+1}$



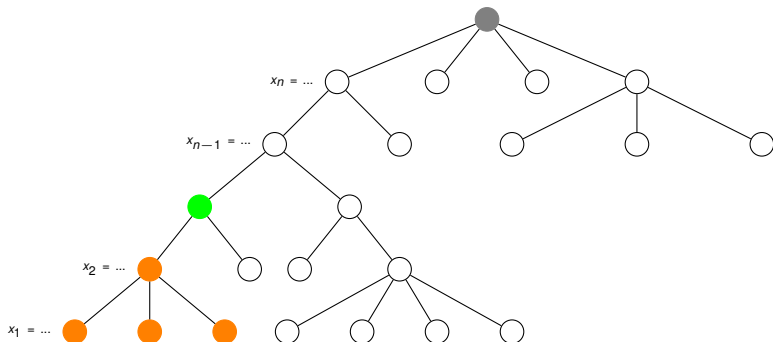
New optimum  $A = \|\mathbf{Bx}\|_2^2$

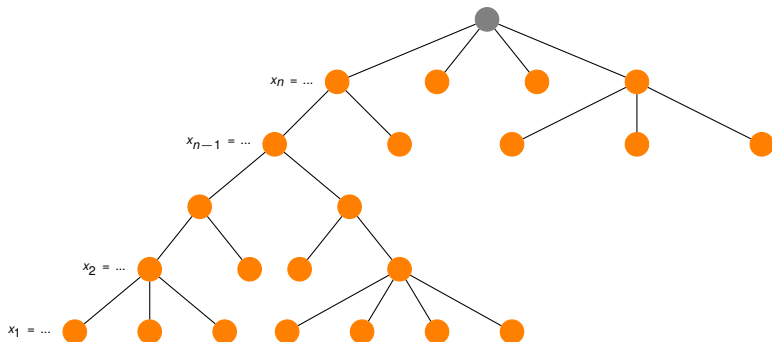


Cut off branch if  $l_j > A$ .

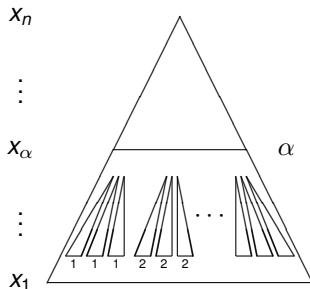








How to find independent subtrees?



# Basic idea

---

**Input:**  $\mathbf{B}$ ,  $A$ ,  $\alpha$ ,  $n$

Compute the Gram-Schmidt decomposition of  $\mathbf{b}_i$

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

---

# Basic idea

---

**Input:**  $\mathbf{B}$ ,  $A$ ,  $\alpha$ ,  $n$

Compute the Gram-Schmidt decomposition of  $\mathbf{b}_i$

CPU: generate  $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$

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

---

# Basic idea

**Input:**  $\mathbf{B}$ ,  $A$ ,  $\alpha$ ,  $n$

Compute the Gram-Schmidt decomposition of  $\mathbf{b}_i$

CPU: generate  $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$

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  $\|\sum_{i=1}^n x_i \mathbf{b}_i\|$  minimal

# Basic idea

**Input:**  $\mathbf{B}$ ,  $A$ ,  $\alpha$ ,  $n$

Compute the Gram-Schmidt decomposition of  $\mathbf{b}_i$

CPU: generate  $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$

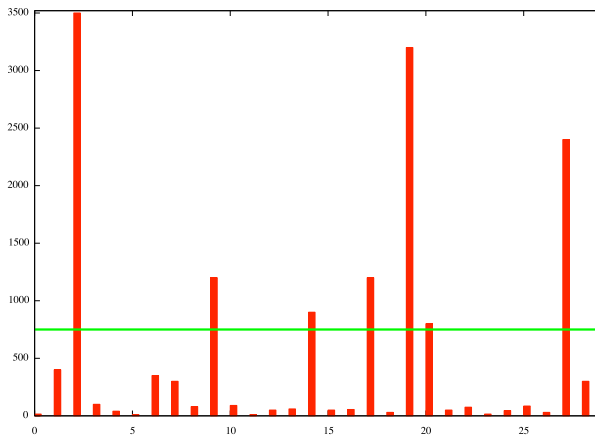
GPU thread: run a sub-enum on  $\mathbf{x}_i$ , if new optimum  $\rightarrow$  store in  $\mathbf{x}$

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

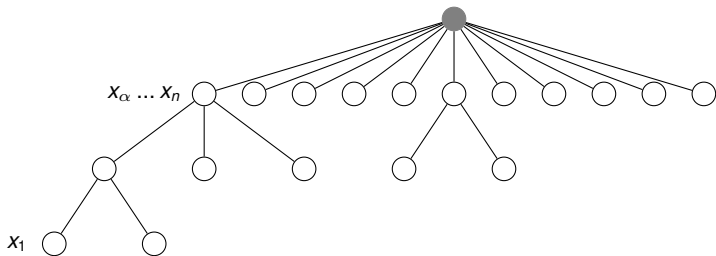
$\Rightarrow$  horrible performance

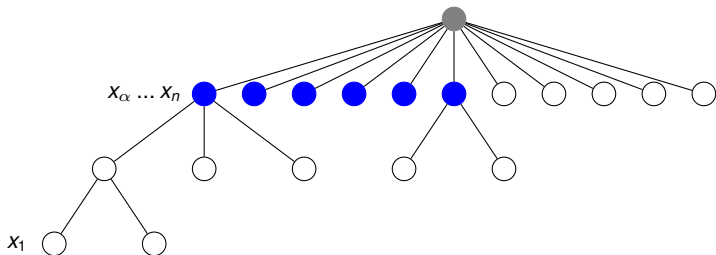


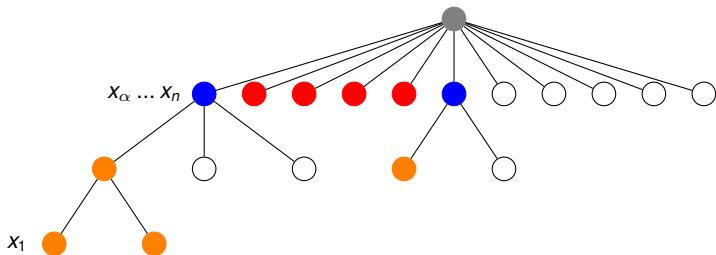
# Thread execution time

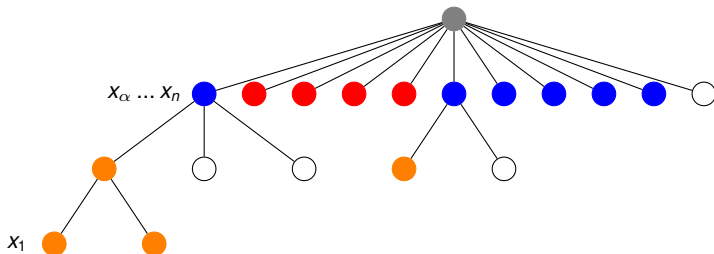


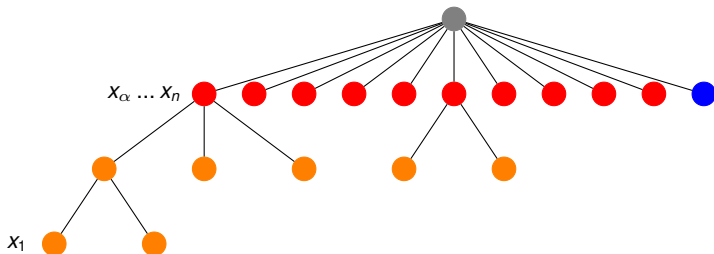












**Input:**  $\mathbf{B}$ ,  $A$ ,  $\alpha$ ,  $n$

Compute the Gram-Schmidt decomposition of  $\mathbf{b}_i$

**while** *true* **do**

CPU: generate some  $\mathbf{x}_i = [0, \dots, 0, x_\alpha, \dots, x_n]$

CPU: Get enum state  $\bar{\mathbf{x}}_i = [\bar{x}_1, \dots, \bar{x}_{\alpha-1}, x_\alpha, \dots, x_n]$

**end**

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

**Input:**  $\mathbf{B}$ ,  $A$ ,  $\alpha$ ,  $n$

Compute the Gram-Schmidt decomposition of  $\mathbf{b}_j$

**while true do**

CPU: generate some  $\mathbf{x}_j = [0, \dots, 0, x_\alpha, \dots, x_n]$

GPU thread:

**while there are  $\mathbf{x}_j$  left do**

Start enum for a certain  $\mathbf{x}_j$  or continue enum for  $\bar{\mathbf{x}}_j$

Stop enum after  $\mathbf{M}$  steps, store the state  $\{l_j, \bar{\mathbf{x}}_j, s_j = \mathbf{M}\}$

**end**

CPU: Get enum state  $\bar{\mathbf{x}}_j = [\bar{x}_1, \dots, \bar{x}_{\alpha-1}, x_\alpha, \dots, x_n]$

**end**

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

Some facts & figures:

- Dimension 50, 100,000 starting vectors → upload & download ~ 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



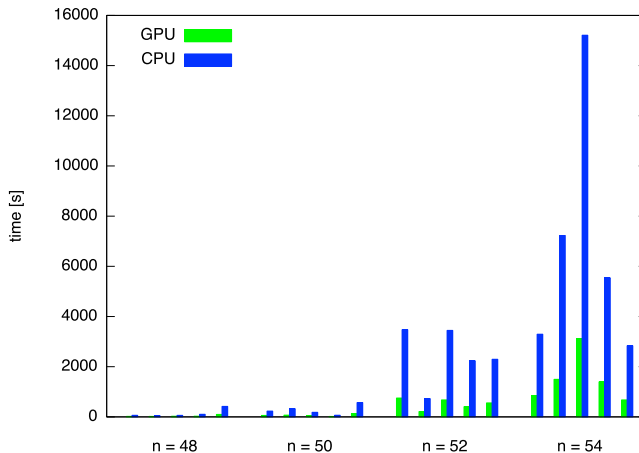
Throughput:

- CPU:  $25 \cdot 10^6$  steps/s
- GPU:  $200 \cdot 10^6$  steps/s

Throughput on GPU depends on:

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

# Results



# Results

n	44	48	52	56	60
fpLLL - ENUM	17.7s	136s	2434s	137489s	-
CUDA - ENUM	11.7s	37.0s	520s	30752s	274268s
	66%	27%	21%	22%	-

Speedup in this model:  $\approx 5$

<http://homes.esat.kuleuven.be/~jhermans/gpuenum/index.html>

## CPU machine

- 4-core CPU
- \$900
- speedup 4
- total cost  $91 \cdot t$

## GPU machine

- 4-core CPU + 2 x GTX295 GPU
- \$2200
- speedup 24
- total cost  $225 \cdot t$

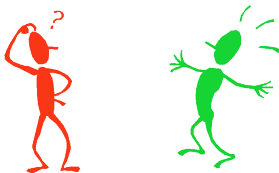
Speedup in this model: **2.4**

- Generalize ideas (not specific for gpu's... clusters?)
- Use full power of CPU (now: idle during gpu-time)
- Test on NVIDIA Fermi cards



(... and we will climb the highest mountains)

# Finally...



Questions?

- GTX 280 can make 13 times more fmadds than Core2Extreme CPU
- using CUDA profiler
- diverging branches, warp serializations, ...
- lots of ALUs left idle