



HPC center of the Netherlands



Augmented Reality Technology

Better Than All the Rest: Finding Maximum-Performance GPU Kernels Using Auto-Tuning

GPU Technology Conference

April 7, 2016

Cedric Nugteren

How to find the best flight ticket?



Which airline?



Fly a day earlier or later?

April 2016 < Today >

Mon	Tue	Wed	Thu	Fri	Sat	Sun
28	29	30	31	1 Apr	2	3
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	1 May
2	3	4	5	6	7	8



Which route?

Solution:

- Evaluate all combinations manually
- ... or use a flight comparison website

How to find the best flight ticket?



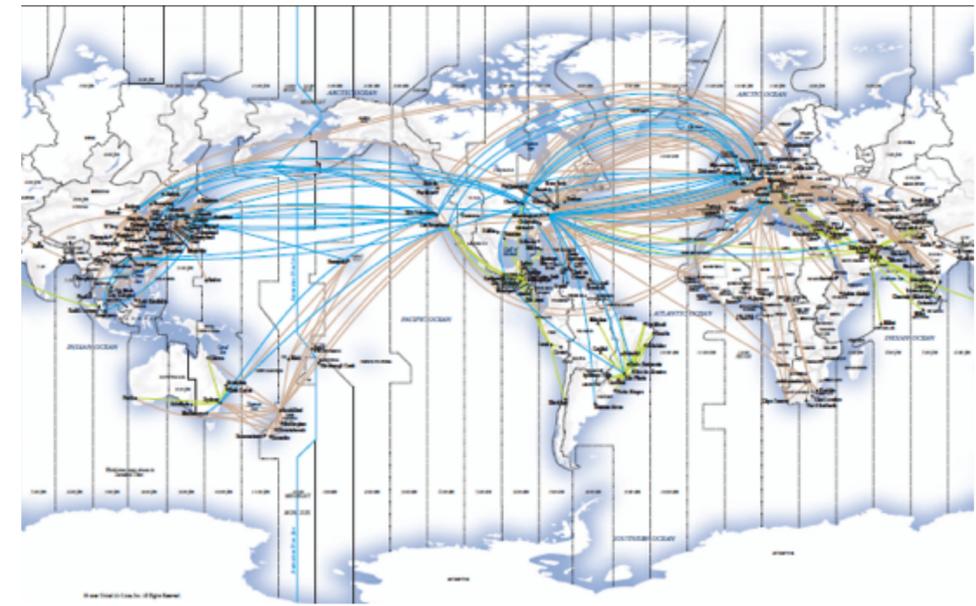
Thread block sizes?



Cache in shared memory?

April 2016

Mon	Tue	Wed	Thu	Fri	Sat	Sun
28	29	30	31	1 Apr	2	3
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	1 May
2	3	4	5	6	7	8



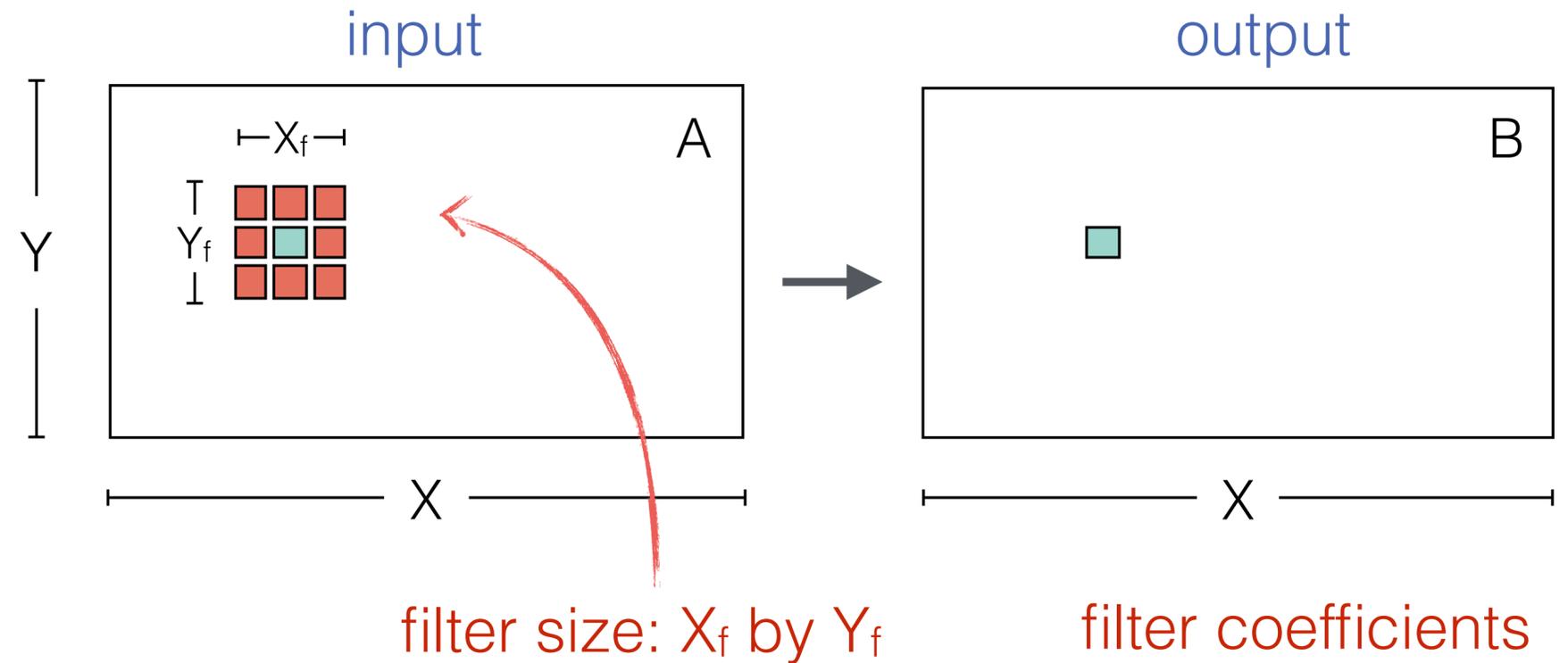
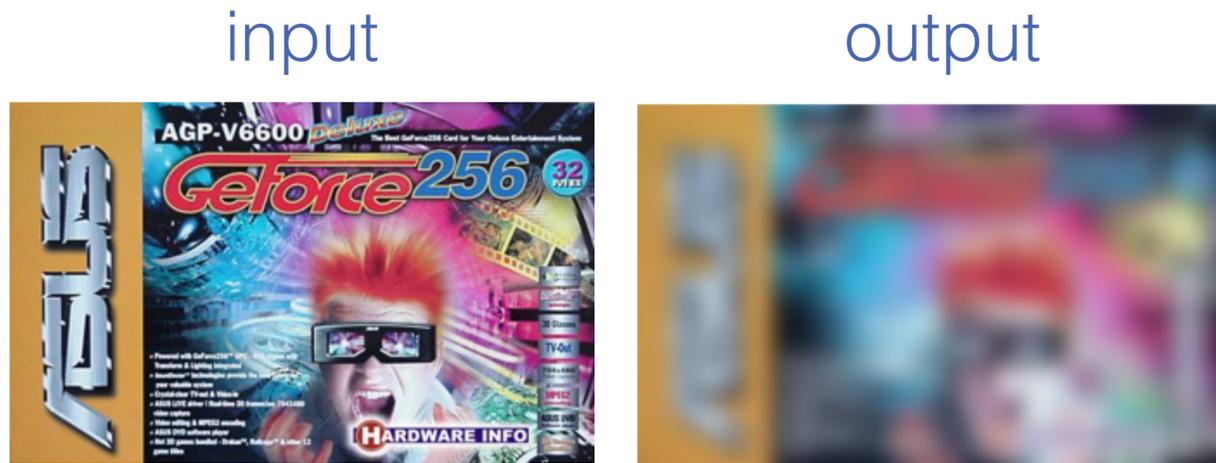
More work per thread or less?

Solution:

- Evaluate all combinations manually
- ... or use a **light CLTune auto-tuner**

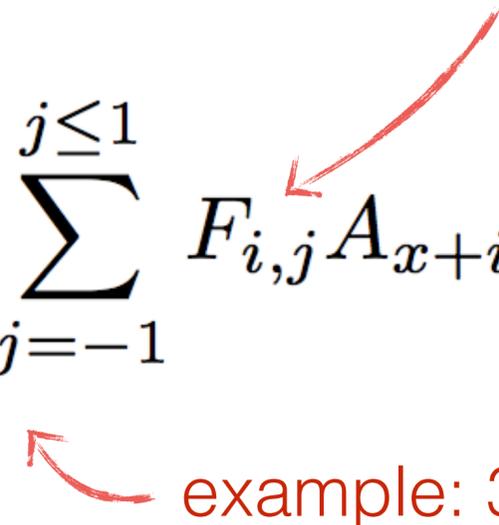
Example: convolution

Example: blur filter



Targets:

- GPUs (CUDA & OpenCL)
- Multi-core CPUs
- Other OpenCL-capable devices

$$B_{x,y} = w \cdot \sum_{i=-1}^{i=1} \sum_{j=-1}^{j=1} F_{i,j} A_{x+i,y+j}$$


example: 3 by 3 filter

OpenCL 2D convolution

each thread: one output pixel

Thread coarsening (2D)?

$$B_{x,y} = w \cdot \sum_{i=-1}^{i \leq 1} \sum_{j=-1}^{j \leq 1} F_{i,j} A_{x+i,y+j}$$

double for-loop

Unroll loops?

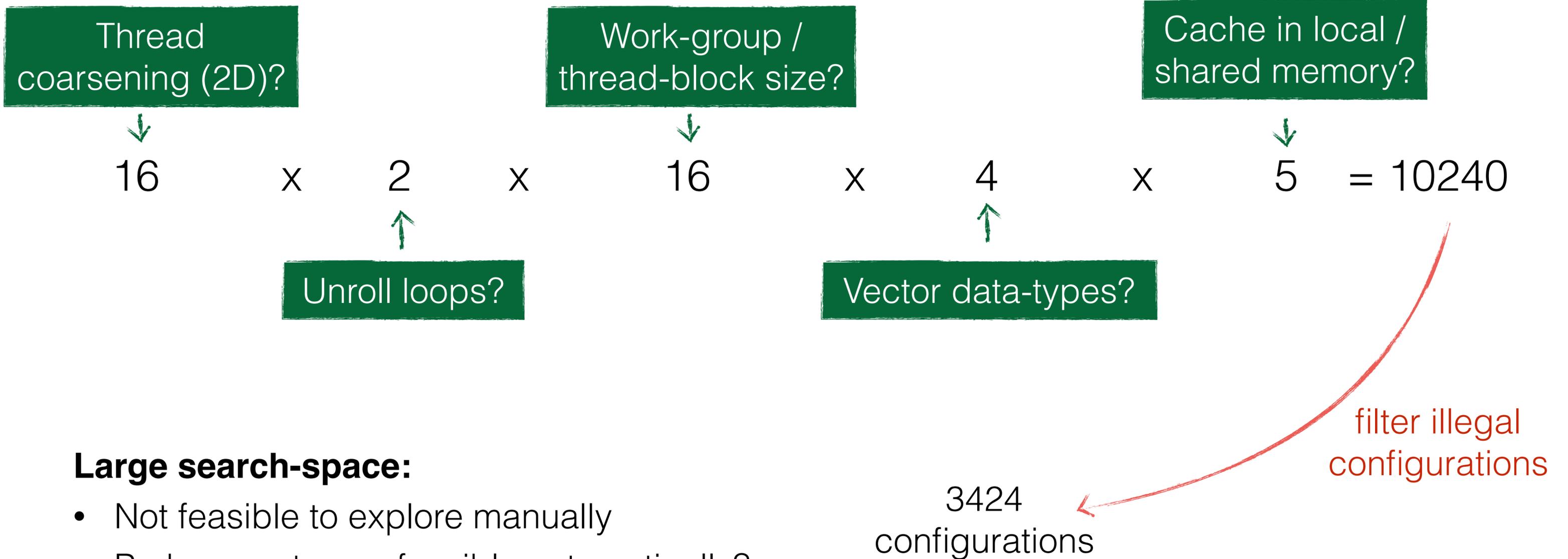
```
1
2 #define HFS (3) // Half filter size
3 #define FS (HFS+HFS+1) // Filter size
4
5 __kernel void conv_reference(const int size_x, const int size_y,
6                             const __global float* src,
7                             __constant float* coeff,
8                             __global float* dest) {
9
10  const int tid_x = get_global_id(0);
11  const int tid_y = get_global_id(1);
12
13  float acc = 0.0f;
14
15  // Loops over the neighbourhood
16  for (int fx=-HFS; fx<=HFS; ++fx) {
17    for (int fy=-HFS; fy<=HFS; ++fy) {
18      const int index_x = tid_x + HFS + fx;
19      const int index_y = tid_y + HFS + fy;
20
21      // Performs the accumulation
22      float coefficient = coeff[(fy+HFS)*FS + (fx+HFS)];
23      acc += coefficient * src[index_y*size_x + index_x];
24    }
25  }
26
27  // Stores the result
28  dest[tid_y*size_x + tid_x] = acc;
29 }
```

work-group / thread-block size?

Vector data-types?

Cache in local / shared memory?

Search-space explosion



Large search-space:

- Not feasible to explore manually
- Perhaps not even feasible automatically?

Why do we need an auto-tuner?

Large search-space:

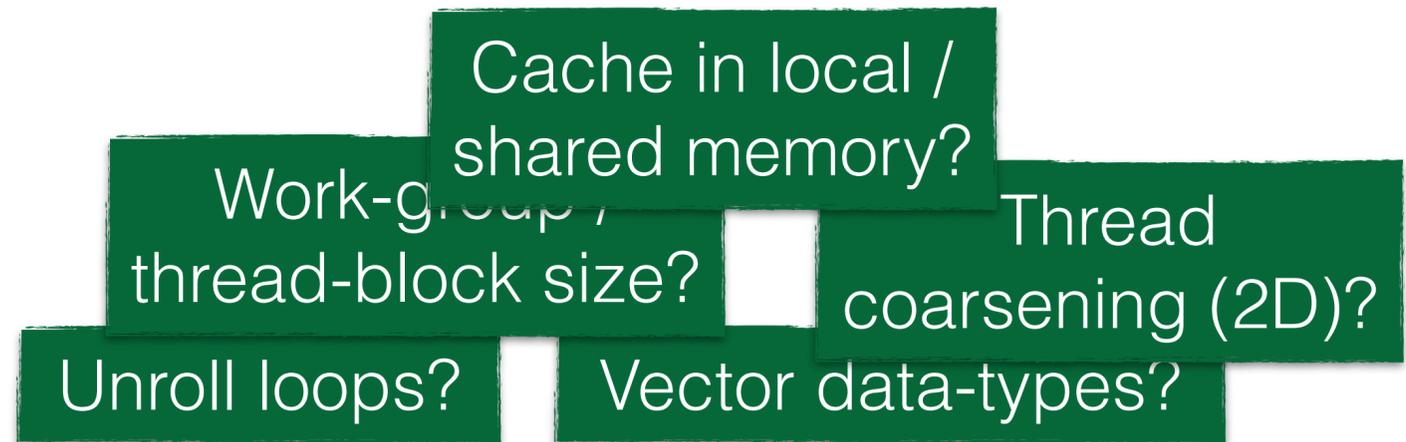
- Not feasible to explore manually
- Perhaps not even feasible automatically?

Wide variety of devices:

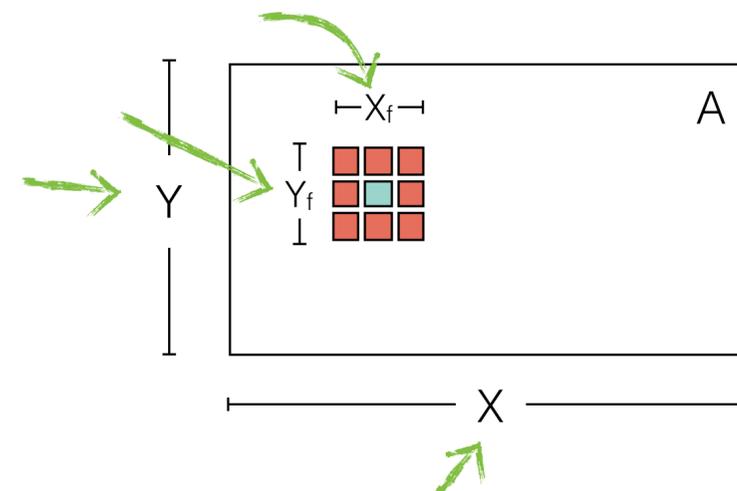
- Different optimal kernels
- Even from the same vendor

User-parameter dependent:

- Examples: matrix sizes, image size, filter sizes, etc.



vendor and device name	architecture	peak GFLOPS	peak GB/s	GFLOPS per GB/s
NVIDIA Tesla K40m	Kepler	4291	288	14.9
NVIDIA GeForce GTX480	Fermi	1345	177	7.6
AMD Radeon HD7970	Tahiti	4368	288	15.1
Intel Iris 5100	Iris	832	26	32.5



Why

Large search-s

- Not feasible to
- Perhaps not e

Wide variety of

- Different opti
- Even from the

User-paramete

- Examples: m
- image size, fi

Repository: CNugteren / CLTune

Statistics: 213 commits, 2 branches, 15 releases, 1 contributor

Branch: master

File/Folder	Description	Commit Date
cmake/Modules	New custom FindOpenCL.cmake	10 months ago
include	Added experimental support for CUDA kernels	4 months ago
samples	Fixed a warning and error for MSVC	4 months ago
src	Prepared for MSVC support	4 months ago
test	Updated the 'KernelInfo' class to use Catch	3 months ago
.gitignore	Added gitignore for build directory	a year ago
.travis.yml	Made Travis always build pushes to the master branch	6 months ago
CHANGELOG	Updated to version 2.0.0	3 months ago
CMakeLists.txt	Updated to version 2.0.0	3 months ago
LICENSE	Updated license information	a year ago
README.md	Updated the readme	3 months ago

CLTune: Automatic OpenCL kernel tuning

build passing

CLTune is a C++ library which can be used to automatically tune your OpenCL and CUDA kernels. The only thing you'll need to provide is a tuneable kernel and a list of allowed parameters and values.

For example, if you would perform loop unrolling or local memory tiling through a pre-processor define, just remove the

inner?

Thread
rsening (2D)?
-types?

OPS
GB/s
.9
.6
.1
.5

CLTune in action

Example: matrix-vector multiplication ($y = Ax$) with tiling

```
1
2 #include "cltune.h"
3
4 void main() {
5     constexpr auto kSizeM = 2048;
6     constexpr auto kSizeN = 4096;
7
8     std::vector<float> mat_a(kSizeN*kSizeM);
9     std::vector<float> vec_x(kSizeN);
10    std::vector<float> vec_y(kSizeM);
11
12    cltune::Tuner tuner(0, 1); // platform 0, device 1
13
14    // Adds a kernel which supports tiling by a factor TS
15    auto id = tuner.AddKernel({"../tiled.opencl"}, "matvec_tiled", {kSizeM}, {1});
16    tuner.AddParameter(id, "TS", {32, 64, 128, 256});
17    tuner.MulLocalSize(id, {"TS"});
18
19    // Sets the function's arguments
20    tuner.AddArgumentScalar(static_cast<int>(kSizeM));
21    tuner.AddArgumentScalar(static_cast<int>(kSizeN));
22    tuner.AddArgumentInput(mat_a);
23    tuner.AddArgumentInput(vec_x);
24    tuner.AddArgumentOutput(vec_y);
25
26    // Runs the tuner
27    tuner.Tune();
28    tuner.PrintToScreen();
29 }
30
```

Example input data

Kernel to tune

Verifies output based on optional reference kernel

Tile size (TS):
4 possible values

CLTune in action

```
1
2 #include "cltune.h"
3
4 void main() {
5     constexpr auto kSizeM = 2048;
6     constexpr auto kSizeN = 4096;
7
8     std::vector<float> mat_a(kSizeN*kSizeM);
9     std::vector<float> vec_x(kSizeN);
10    std::vector<float> vec_y(kSizeM);
11
12    cltune::Tuner tuner(0, 1); // platform 0, device 1
13
14    // Adds a kernel which supports tiling by a factor TS
15    auto id = tuner.AddKernel({"../tiled.opencl"}, "matvec_tiled", {kSizeM}, {1});
16    tuner.AddParameter(id, "TS", {32, 64, 128, 256});
17    tuner.MulLocalSize(id, {"TS"});
18
19    // Sets the function's arguments
20    tuner.AddArgumentScalar(static_cast<int>(kSizeM));
21    tuner.AddArgumentScalar(static_cast<int>(kSizeN));
22    tuner.AddArgumentInput(mat_a);
23    tuner.AddArgumentInput(vec_x);
24    tuner.AddArgumentOutput(vec_y);
25
26    // Runs the tuner
27    tuner.Tune();
28    tuner.PrintToScreen();
29 }
30
```

Handles all the host -
device data transfers!

```
1
2 __kernel void matvec_tiled(const int M, const int N,
3                             const __global float* mat_a,
4                             const __global float* vec_x,
5                             __global float* vec_y) {
6     const int i = get_global_id(0); // From 0 to M - 1
7     const int tile_id = get_local_id(0); // From 0 to TS - 1
8
9     // Initializes the local/shared memory
10    __local float tile_x[TS];
11
12    // Loops over the tiles
13    float result = 0.0f;
14    for (int t=0; t<N/TS; ++t) {
15
16        // Loads a tile of the vector x into the local memory
17        tile_x[tile_id] = vec_x[t*TS + tile_id];
18        barrier(CLK_LOCAL_MEM_FENCE);
19
20        // Computes the partial result
21        for (int j=0; j<TS; ++j) {
22            result += mat_a[(t*TS + j)*M + i] * tile_x[j];
23        }
24        barrier(CLK_LOCAL_MEM_FENCE);
25    }
26
27    // Stores the result
28    vec_y[i] = result;
29 }
```

Tuneable tile-
size TS

CLTune in action

```
1
2 #include "cltune.h"
3
4 void main() {
5     constexpr auto kSizeM = 2048;
6     constexpr auto kSizeN = 4096;
7
8     std::vector<float> mat_a(kSizeM);
9     std::vector<float> vec_x(kSizeN);
10    std::vector<float> vec_y(kSizeN);
11
12    cltune::Tuner tuner(0, 1); //
13
14    // Adds a kernel which supports
15    auto id = tuner.AddKernel({"matvec_tiled"});
16    tuner.AddParameter(id, "TS", 32);
17    tuner.MulLocalSize(id, {"TS"});
18
19    // Sets the function's arguments
20    tuner.AddArgumentScalar("matvec_tiled", 10);
21    tuner.AddArgumentScalar("matvec_tiled", 9);
22    tuner.AddArgumentInput(mat_a);
23    tuner.AddArgumentInput(vec_x);
24    tuner.AddArgumentOutput(vec_y);
25
26    // Runs the tuner
27    tuner.Tune();
28    tuner.PrintToScreen();
29 }
30
```

```
[=====] Initializing on platform 0 device 2
[=====] Device name: 'AMD Radeon R9 M370X Compute Engine' (OpenCL 1.2 )
[-----] Testing reference matvec_reference
[ RUN ] Running matvec_reference
[ OK ] Completed matvec_reference (11 ms) - 1 out of 1
[-----] Testing kernel matvec_tiled
[ RUN ] Running matvec_tiled
[ OK ] Completed matvec_tiled (10 ms) - 1 out of 4
[ RUN ] Running matvec_tiled
[ OK ] Completed matvec_tiled (9 ms) - 2 out of 4
[ RUN ] Running matvec_tiled
[ OK ] Completed matvec_tiled (8 ms) - 3 out of 4
[ RUN ] Running matvec_tiled
[ OK ] Completed matvec_tiled (8 ms) - 4 out of 4
[-----] Printing results to stdout
[ RESULT ] matvec_tiled; 10 ms; TS 32;
[ RESULT ] matvec_tiled; 9 ms; TS 64;
[ RESULT ] matvec_tiled; 8 ms; TS 128;
[ RESULT ] matvec_tiled; 8 ms; TS 256;
[-----] Printing best result to stdout
[ BEST ] matvec_tiled; 8 ms; TS 128;
[-----] Printing results to file in JSON format
[=====] End of the tuning process
```

```
... M, const int N,
global float* mat_a,
global float* vec_x,
float* vec_y) {
// From 0 to M - 1
0); // From 0 to TS - 1
memory
... into the local memory
tile_id];
... i] * tile_x[j];
```

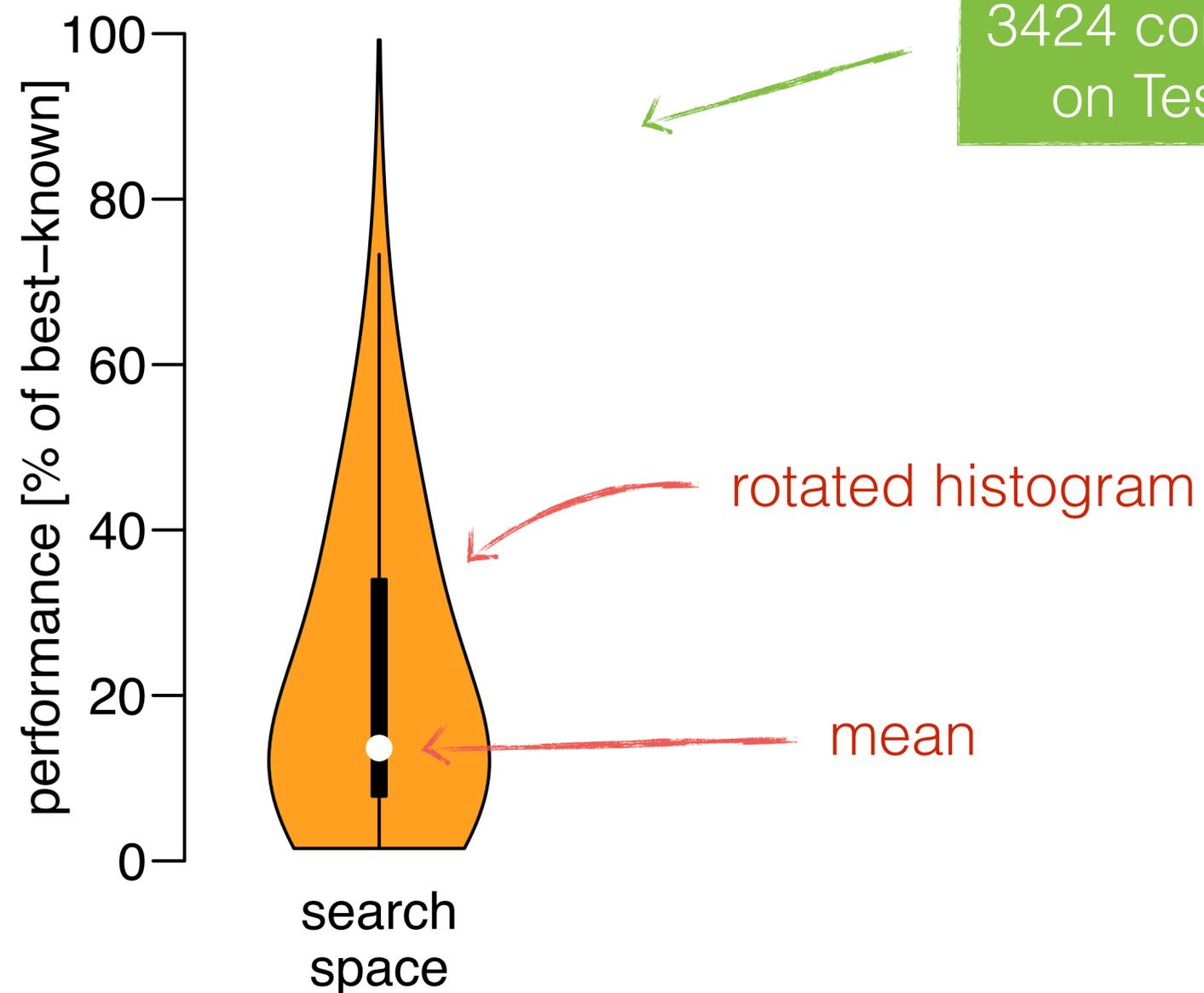
Tuneable tile-size TS

Search strategies

Option 0: Full search

😊 Finds optimal solution

✗ Explores all options



Search strategies

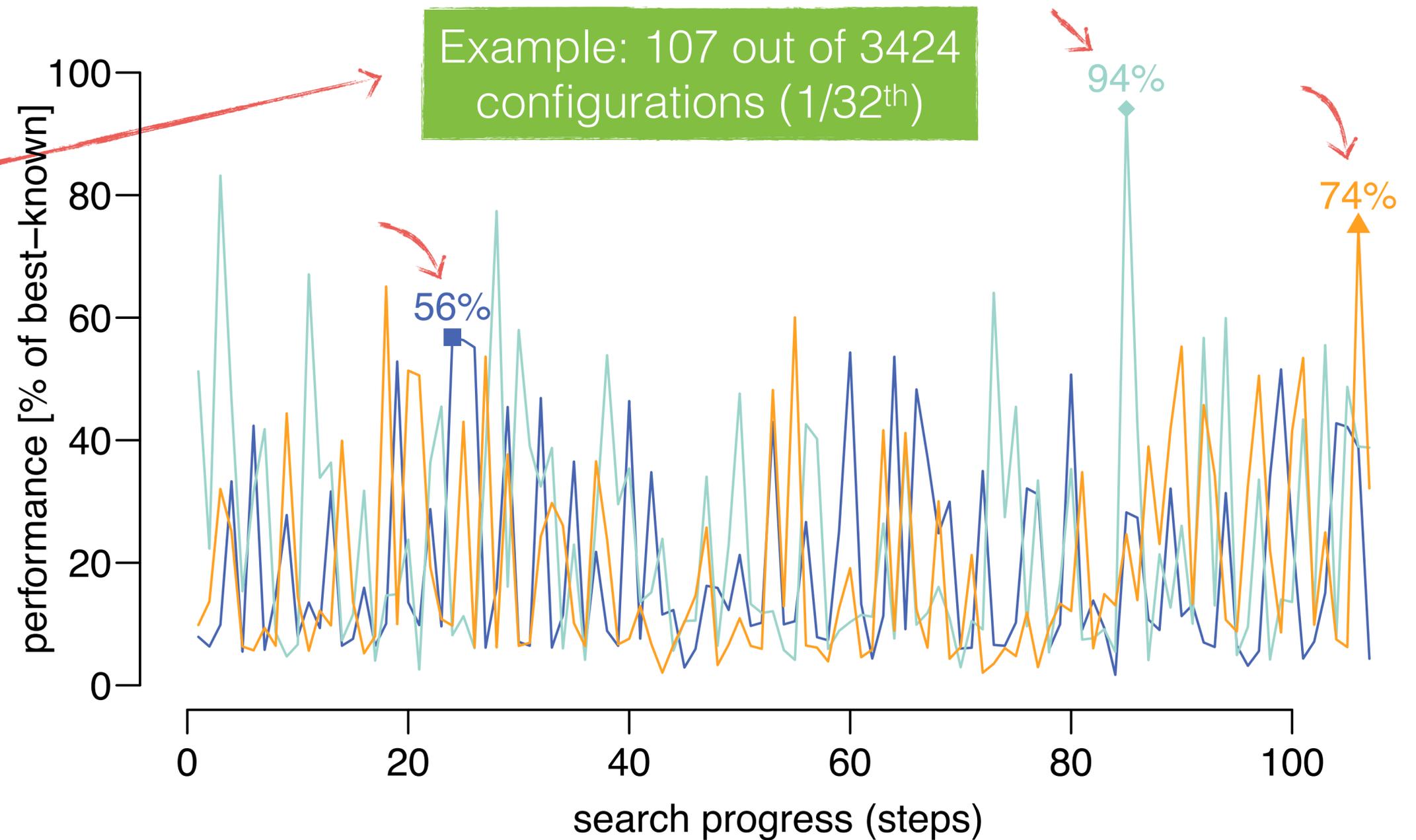
Option 0: Full search

Option 1: Random search

😊 Explores arbitrary fraction

✗ Performance varies

Colours: 3 example runs



Search strategies

Option 0: Full search

Option 1: Random search

Option 2: Simulated annealing

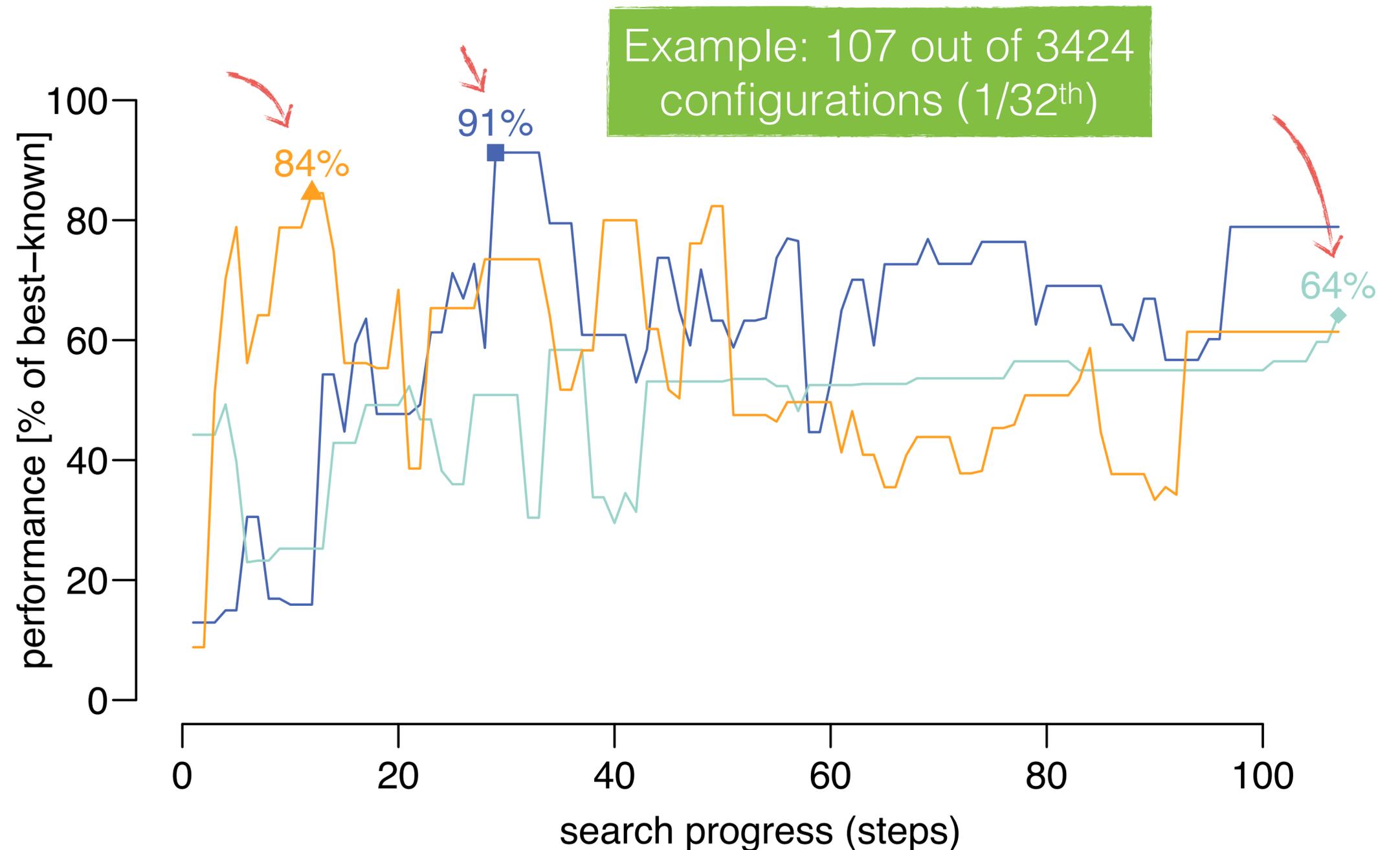
😊 Explores arbitrary fraction

✗ Performance varies

✗ Meta-parameter

✗ Local optima

Colours: 3 example runs



Search strategies

Option 0: Full search

Option 1: Random search

Option 2: Simulated annealing

Option 3: Particle swarm optim

😊 Explores arbitrary fraction

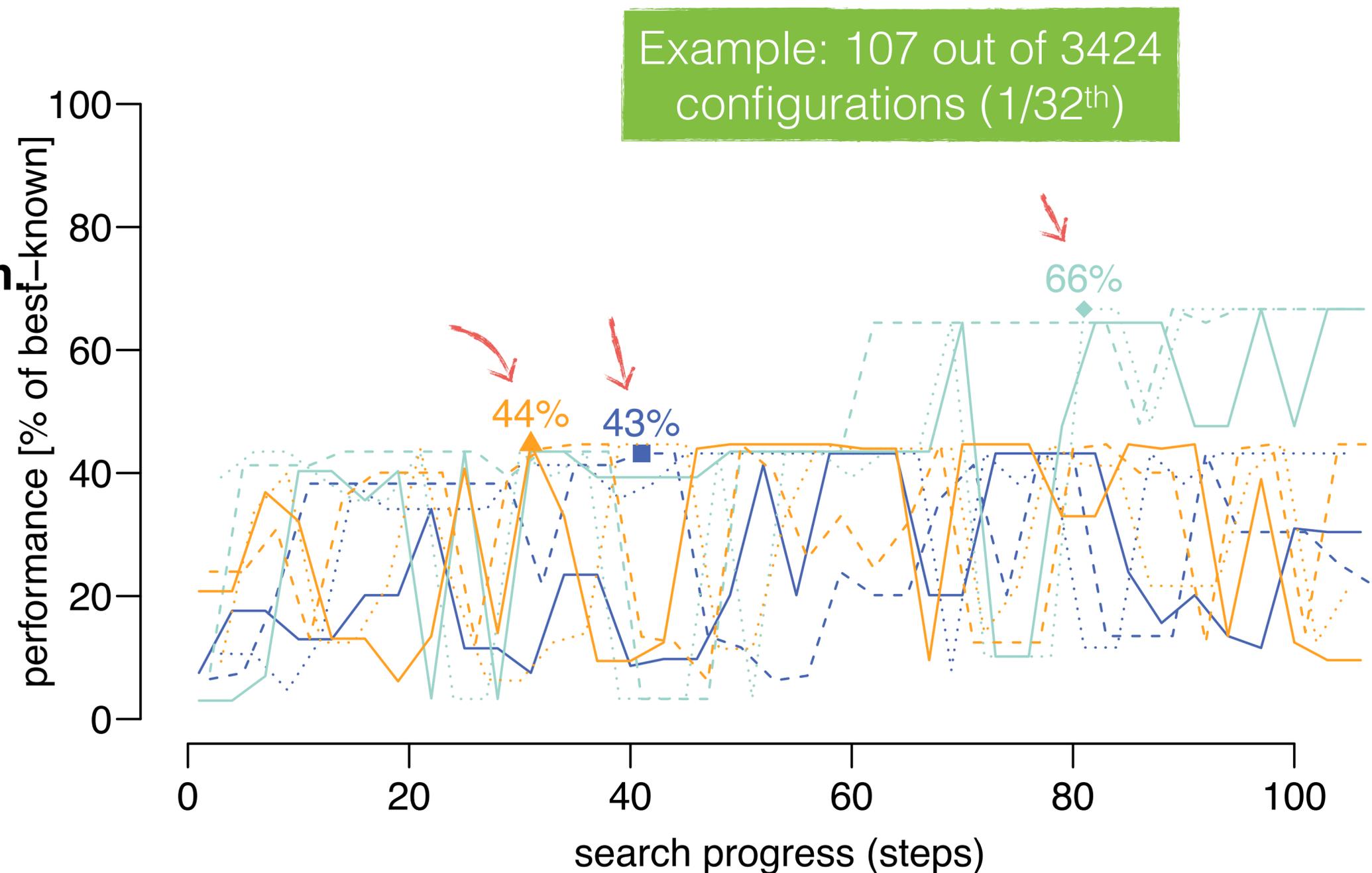
✗ Performance varies

✗ Meta-parameter

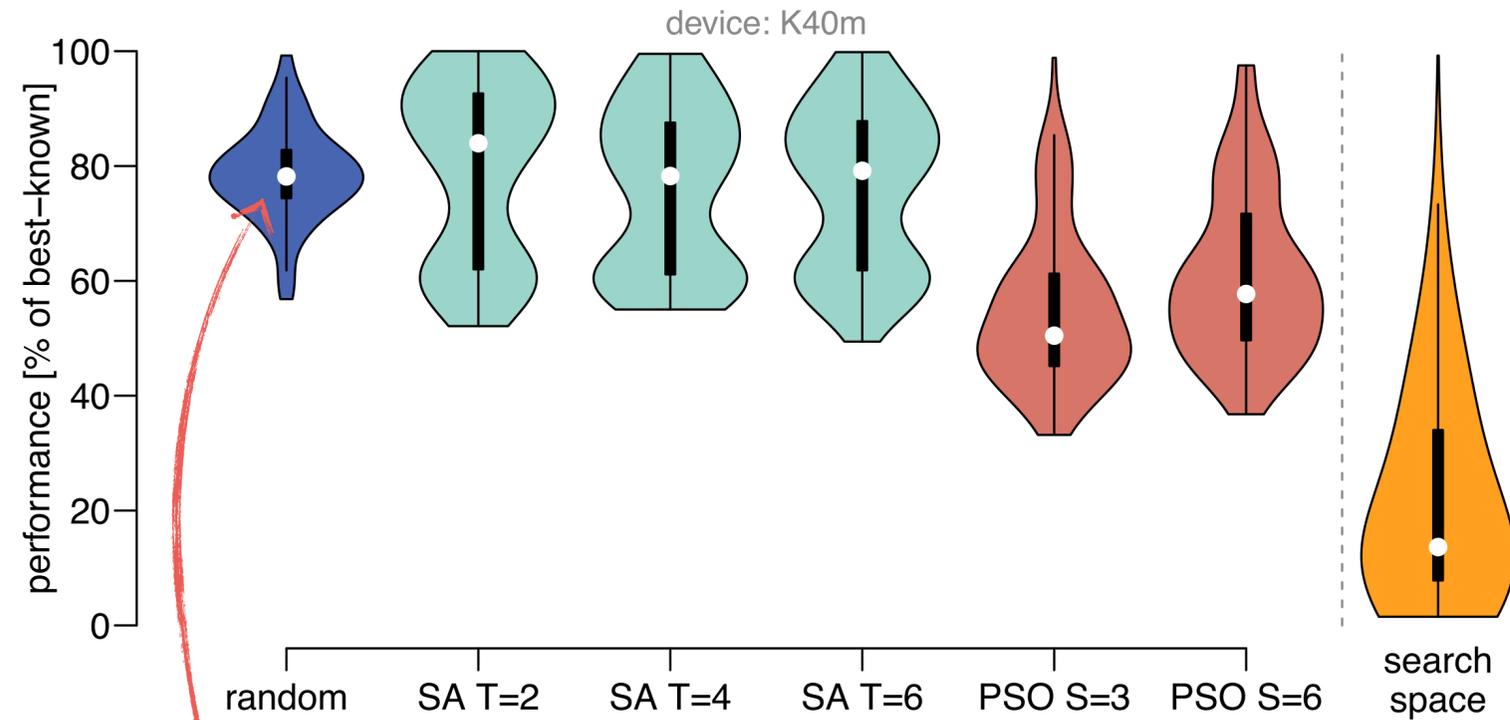
✗ Local optima

Colours: 3 example runs

Line-types: 3 swarms



Search strategies evaluation



average best result
of 128 searches

meta-parameters
for SA and PSO

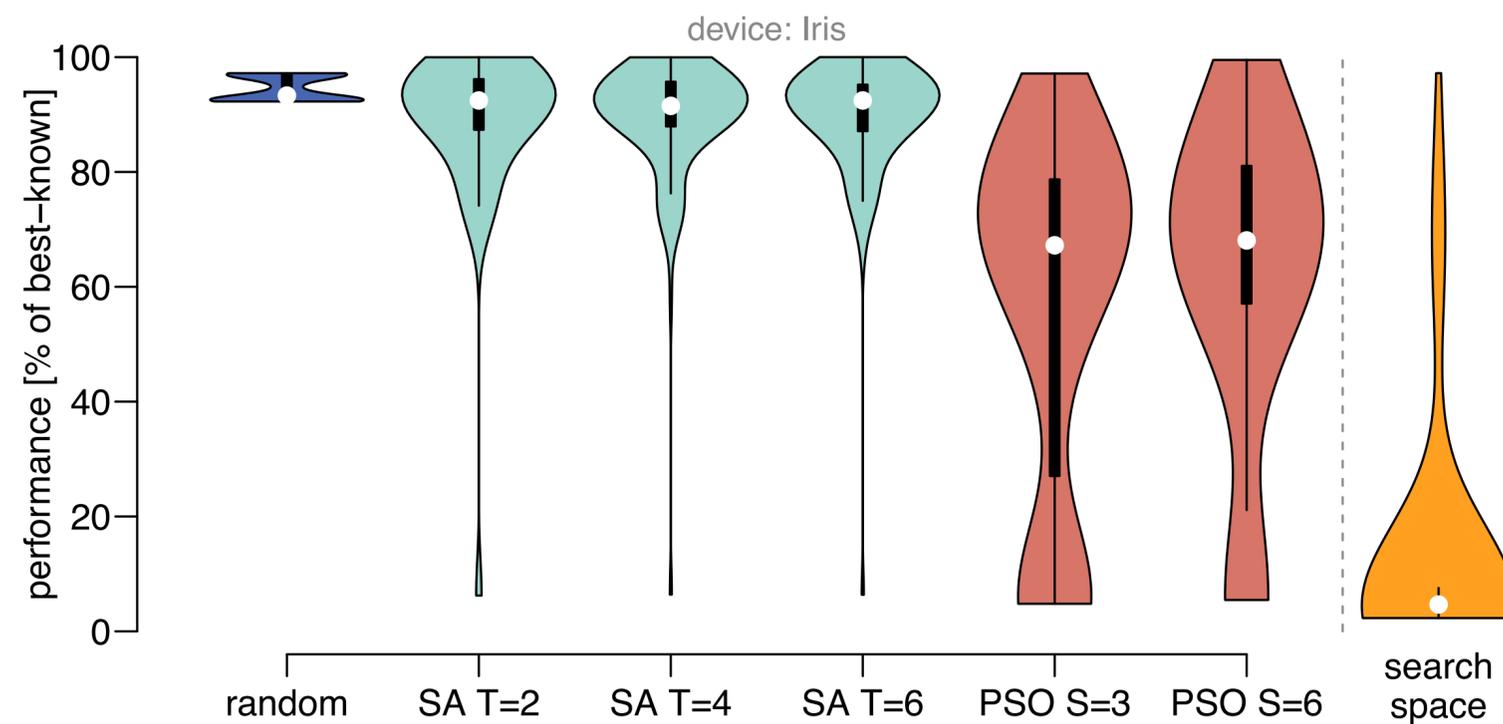
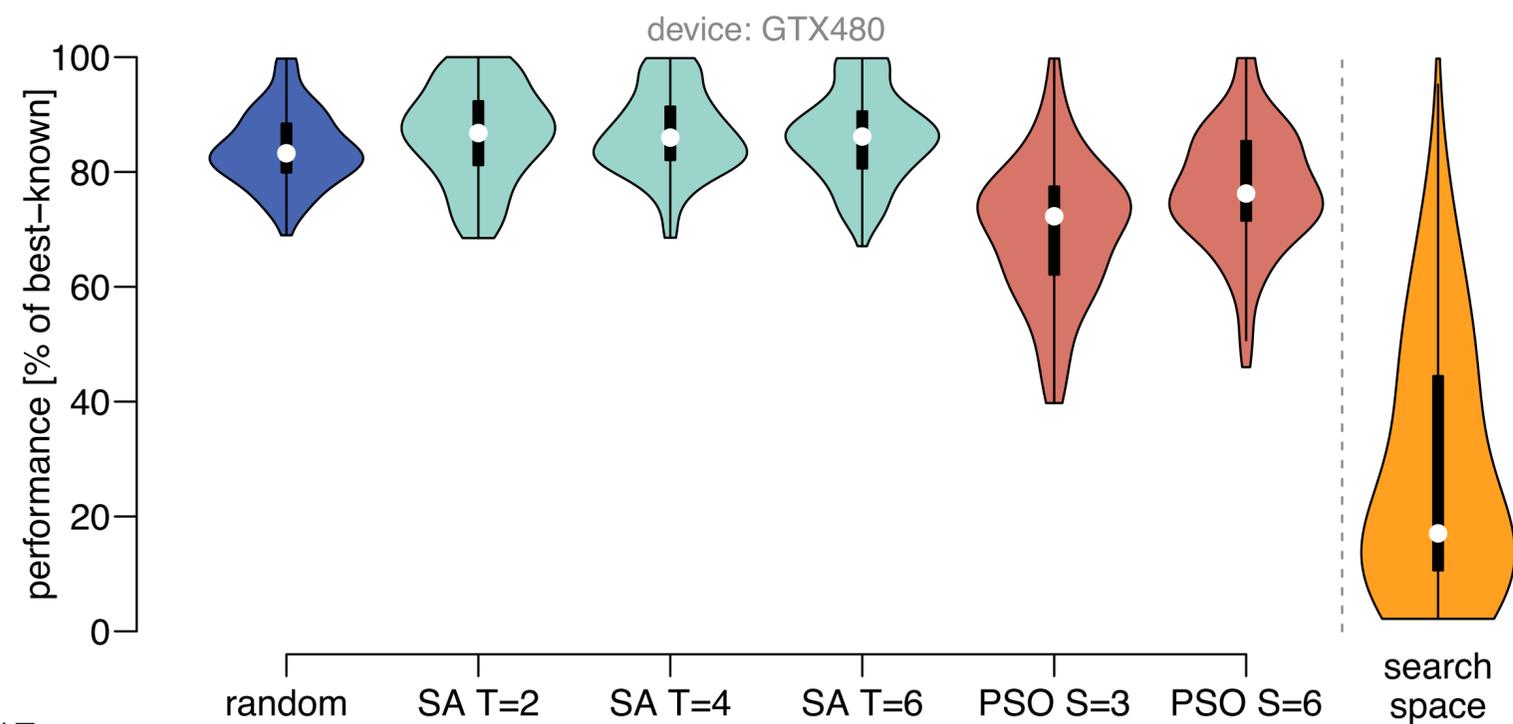
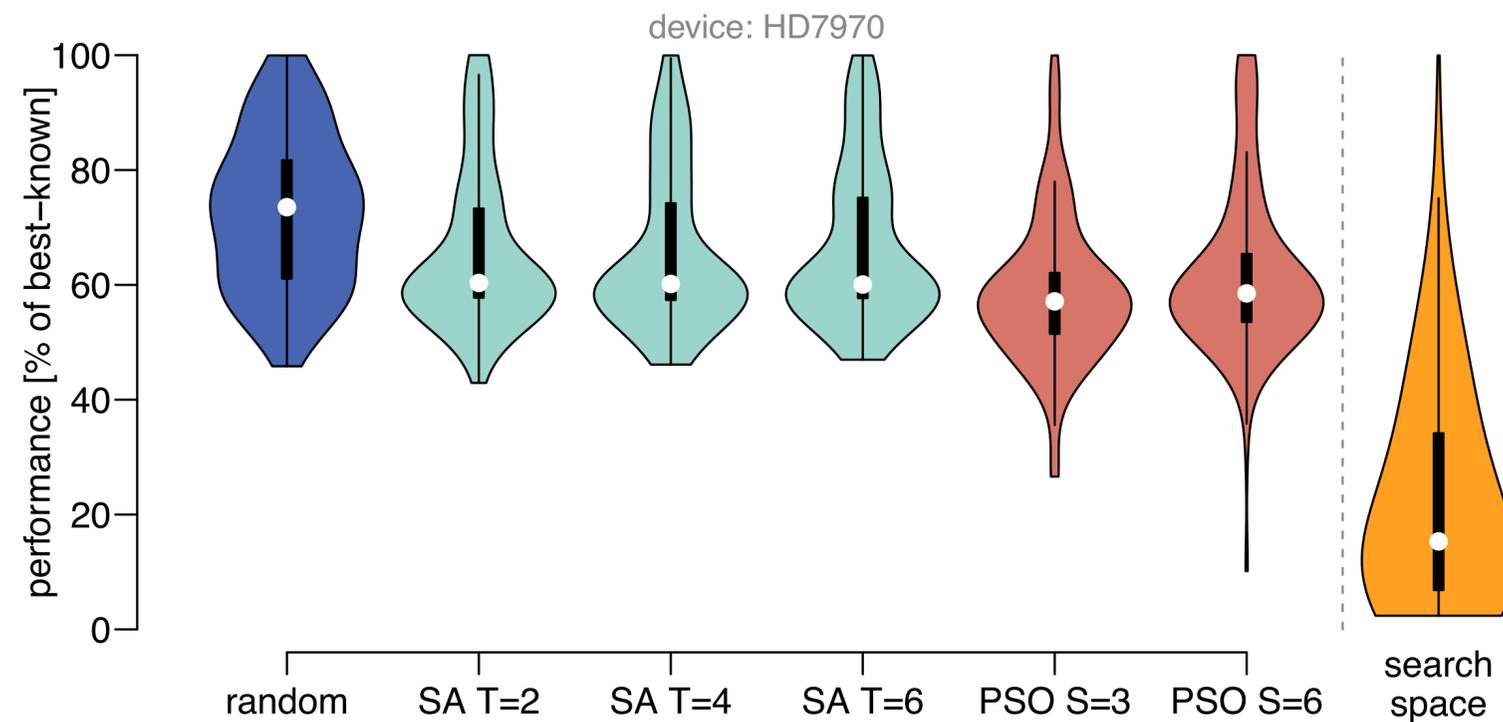
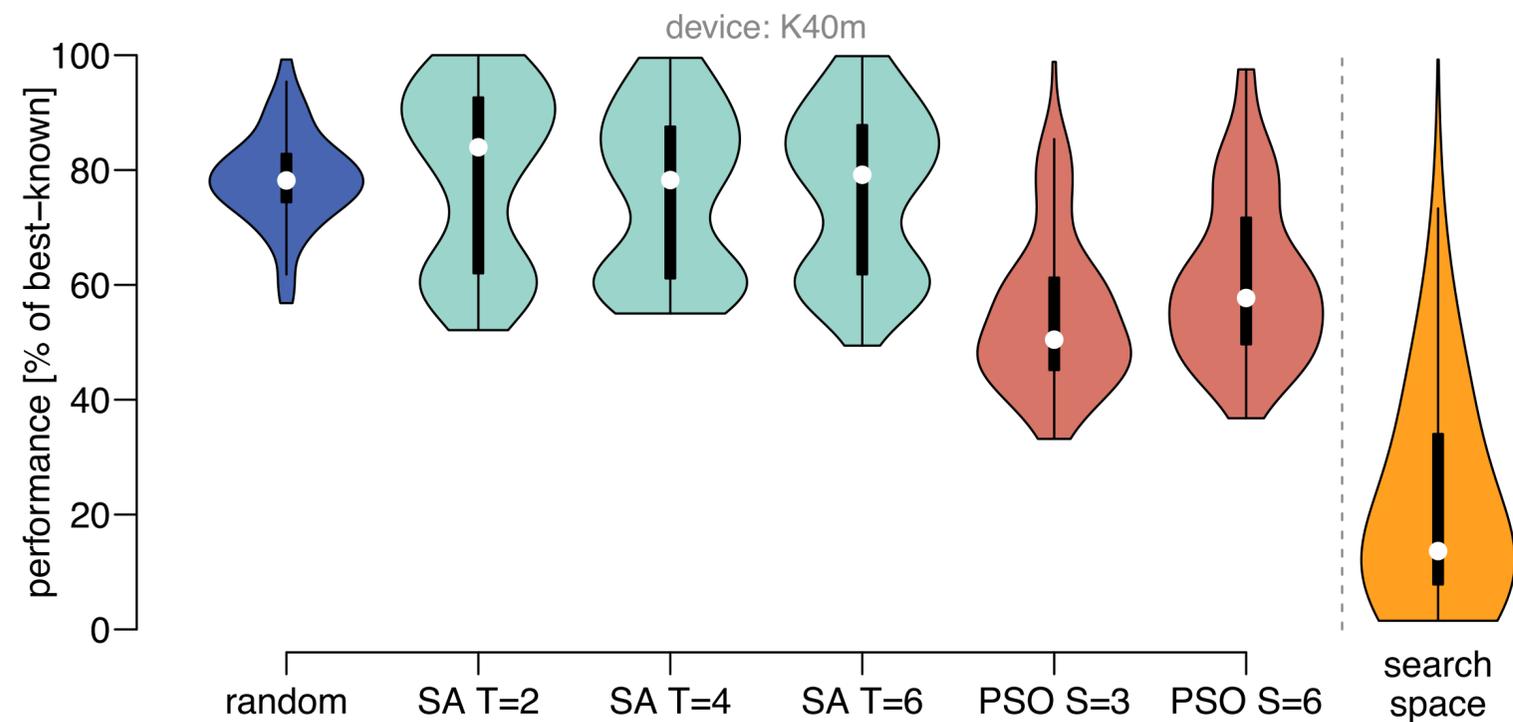
Conclusions:

- PSO performs poorly
- SA perform good, but not much better than random search

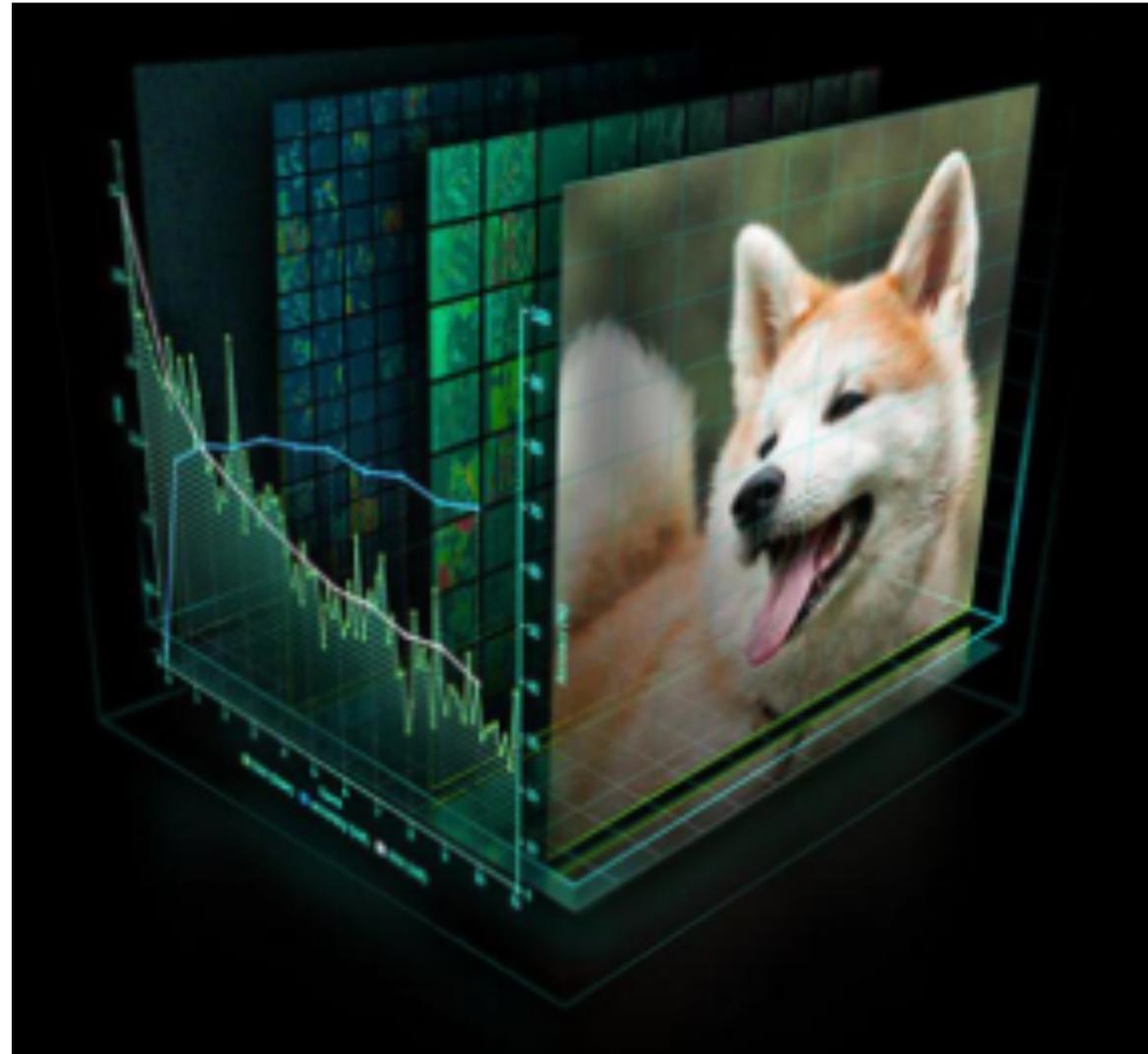
Each search: 107 out of 3424
configurations ($1/32^{\text{th}}$)

[1]: C. Nugteren and V. Codreanu. CLTune: A Generic Auto-Tuner for OpenCL Kernels. IEEE MCSoc '15.

Search strategies evaluation



Solution: machine learning?



Machine learning an auto-tuner

Evaluate subset of all configurations
(e.g. 100 out of 3424)



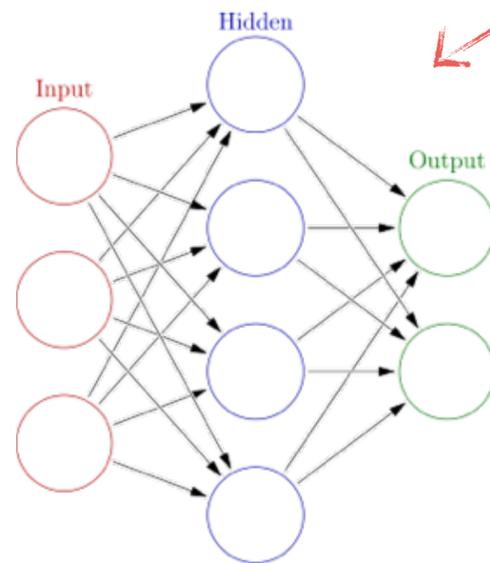
Train model with the examples



Predict execution time for all other configurations



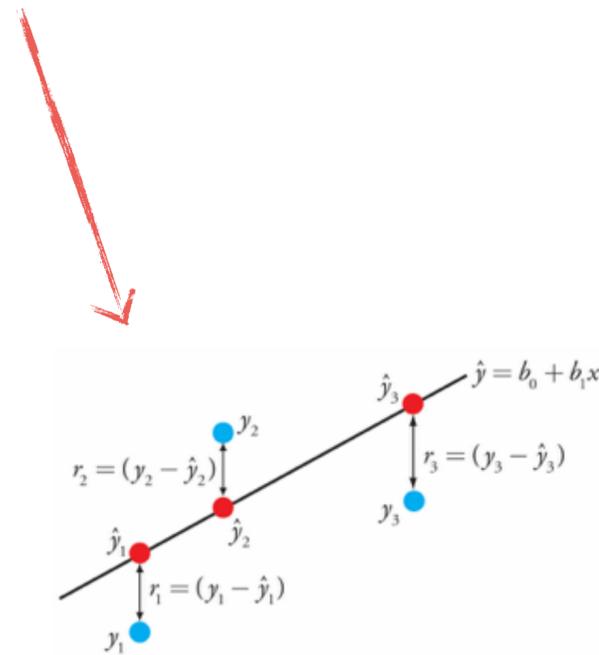
(take best 10 and evaluate them on actual hardware)



Neural network
(3-layer fully connected)

input:
parameter
configuration

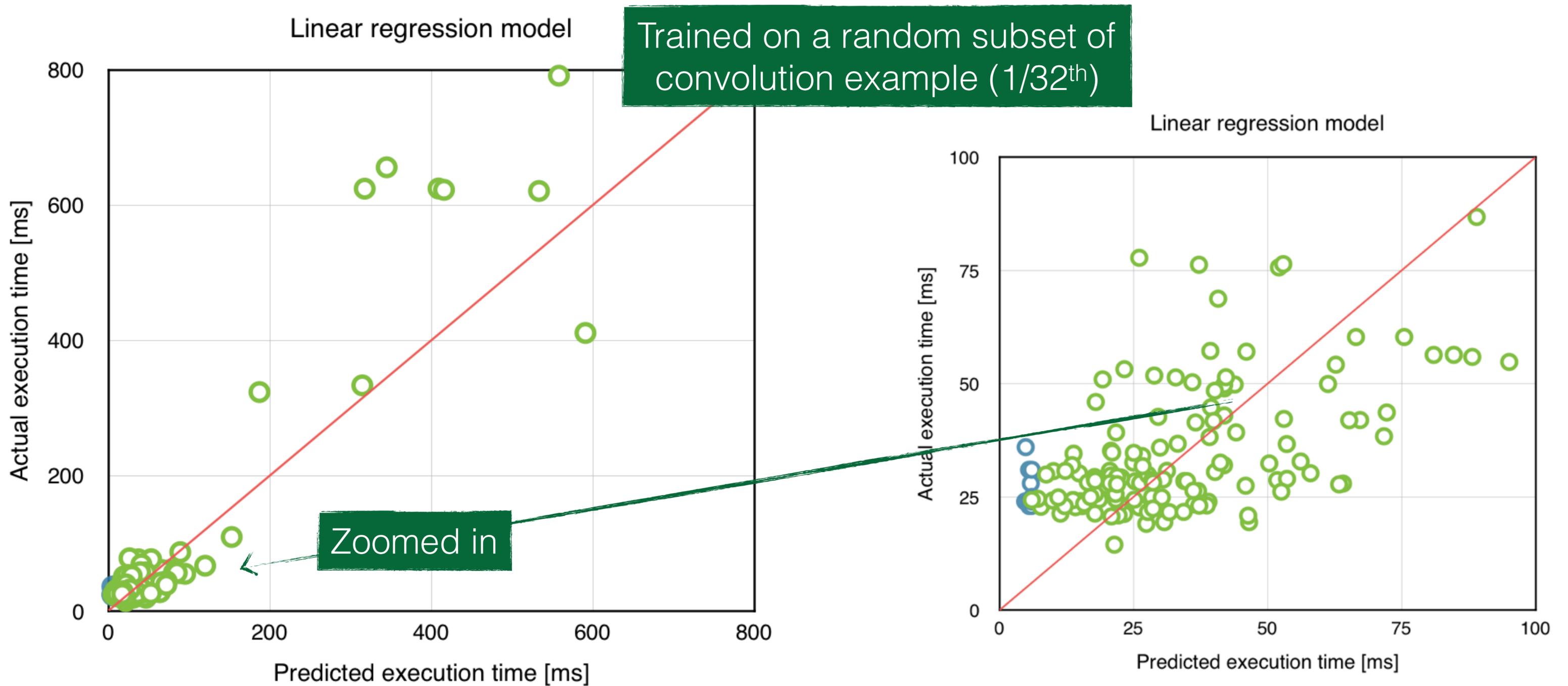
output:
execution time



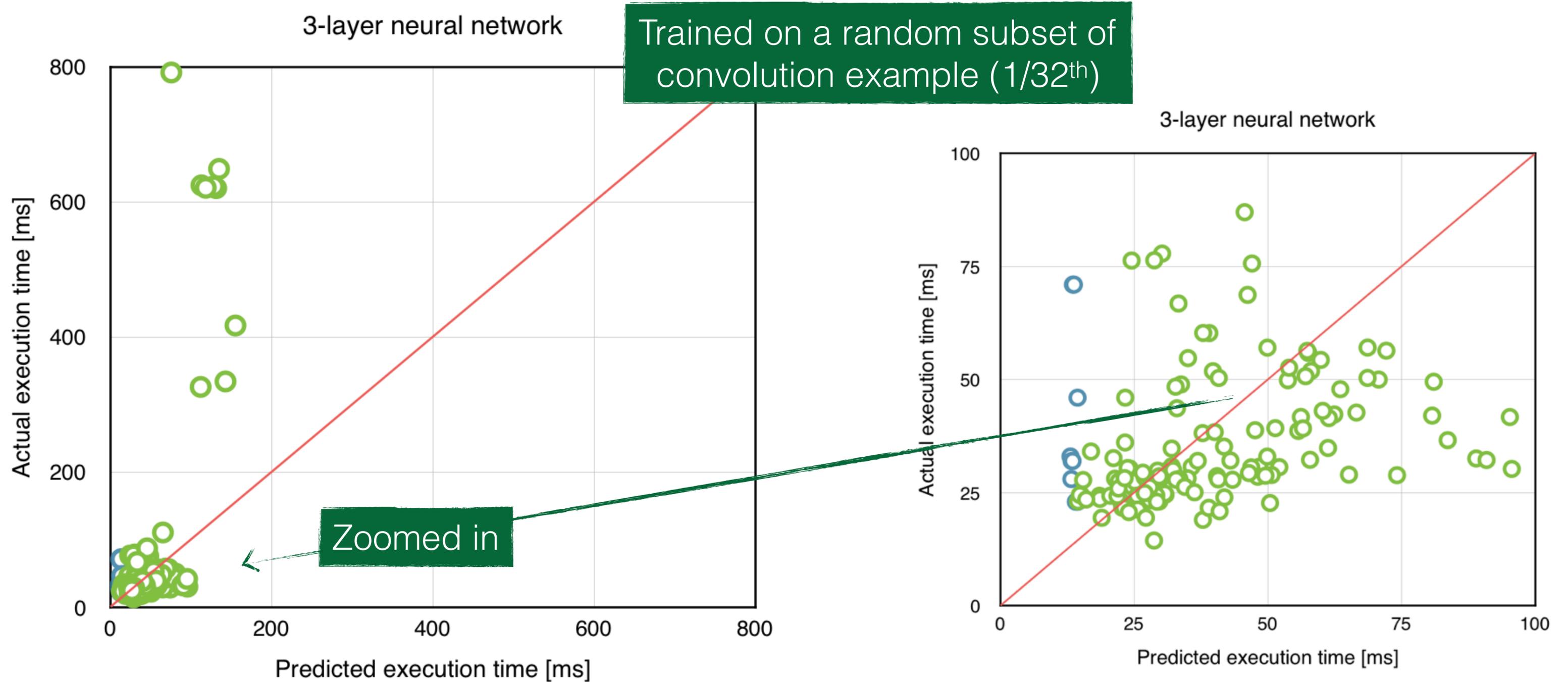
Linear regression

[2]: T.L. Falch and A.C. Elster. Machine Learning Based Auto-tuning for Enhanced OpenCL Performance Portability.

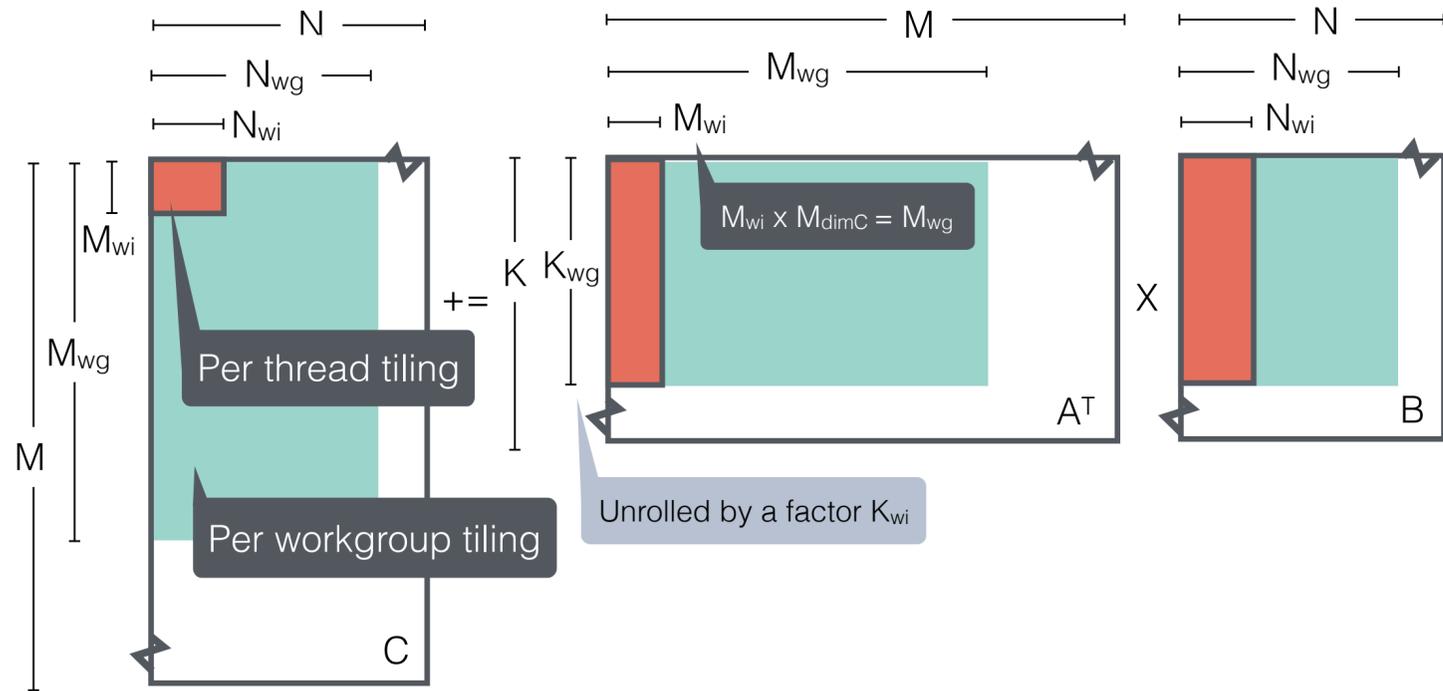
Machine learning an auto-tuner



Machine learning an auto-tuner



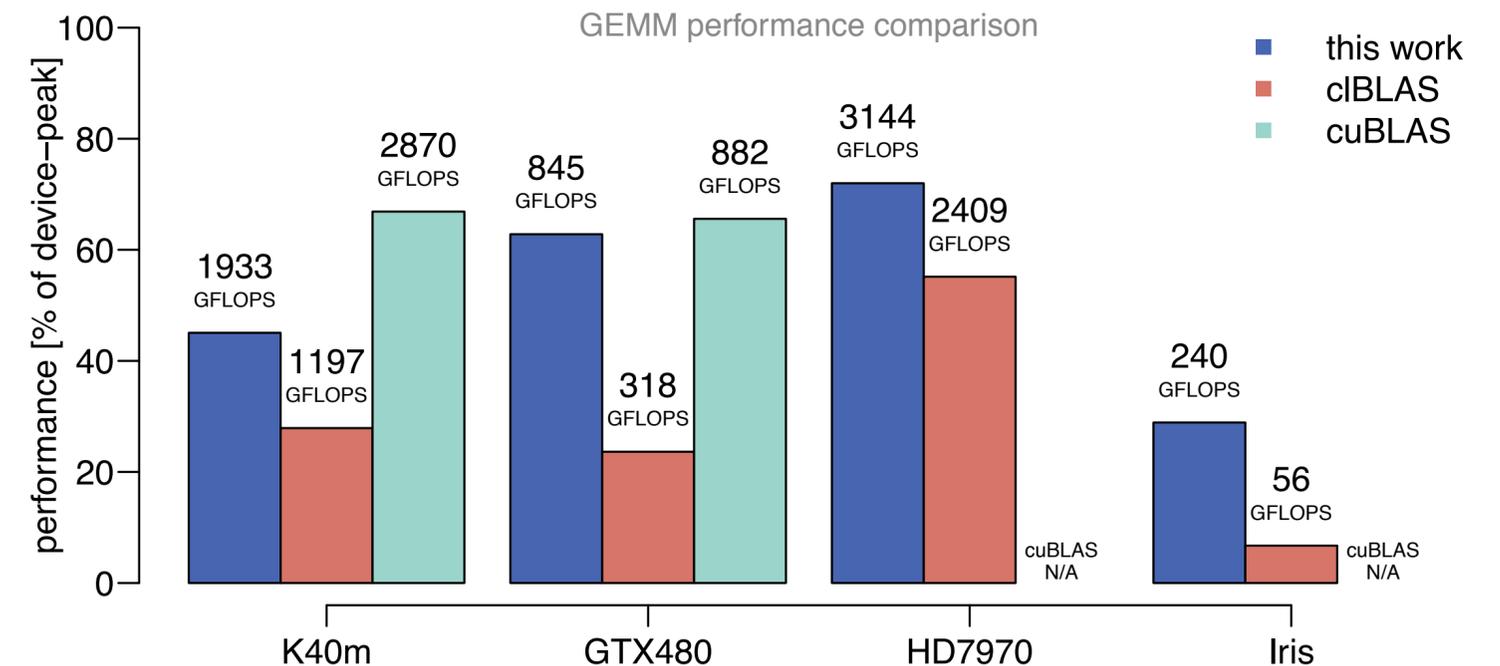
GEMM case-study



Conclusions:

- Different best parameters for different devices
- Performance better than cBLAS, but not as good as assembly-tuned cuBLAS

parameter(s)	allowed values	best parameters per device			
		K40m	GTX480	HD7970	Iris
M_{wg}, N_{wg}, K_{wg}	{16,32,64,128}	128,128,16	64,64,32	128,128,32	64,64,16
M_{dimC}, N_{dimC}	{8,16,32}	16,16	8,16	16,16	8,8
$L\$A, L\B	{yes,no}	yes, yes	yes, yes	yes, yes	yes, yes
M_{dimA}, N_{dimB}	{8,16,32}	32,16	32,32	32,32	8,16
M_{stride}, N_{stride}	{yes,no}	yes, no	yes, no	no, yes	yes, yes
M_{vec}, N_{vec}	{1,2,4,8}	2,1	2,2	4,4	4,4
K_{wi}	{2,8}	8	8	2	8



Tuned OpenCL BLAS — Edit

237 commits 2 branches 6 releases 1 contributor

Branch: master New pull request New file Upload files Find file SSH git@github.com:CNugtere Download ZIP

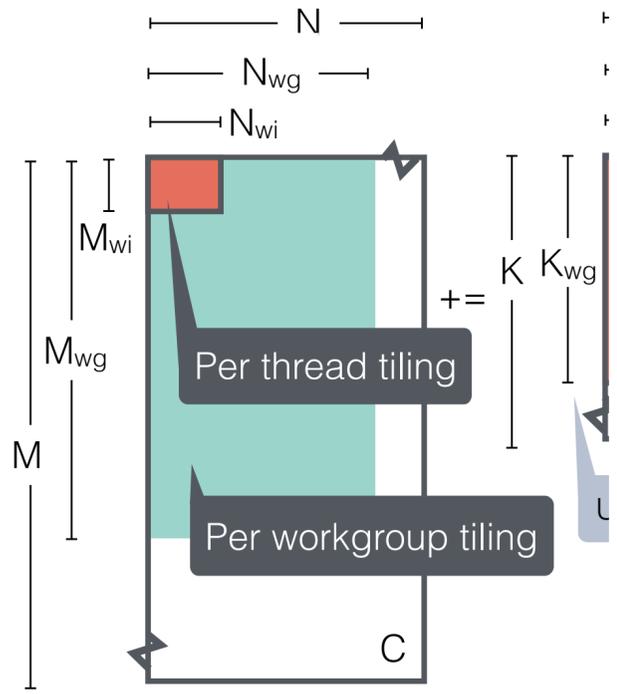
- CNugteren Merge pull request #31 from CNugteren/development Latest commit d190bec 5 hours ago
- cmake/Modules Fixes for compilation under Visual Studio a month ago
 - doc/performance Added performance graphs for Intel Iris and Radeon M370X a day ago
 - include Added tuning results for the newest xGER family kernels a day ago
 - samples Added SGEMM example using the C API 7 months ago
 - scripts Added preliminary support for xHPR2 and xSPR2 routines 7 days ago
 - src Fixed a bug in the GER-family of routines due to incorrect division o... 7 days ago
 - test Made testing against cBLAS in the client binaries truly optional (w... 7 days ago
 - .gitignore Made the tuning database an optional external download a month ago
 - .travis.yml Updated Travis to reflect the changes in the Khronos website 5 hours ago
 - CHANGELOG Updated to version 0.6.0 5 hours ago
 - CMakeLists.txt Updated to version 0.6.0 5 hours ago
 - LICENSE Initial commit of preview version 10 months ago
 - README.md Updated the README file 5 hours ago

README.md

CLBlast: The tuned OpenCL BLAS library

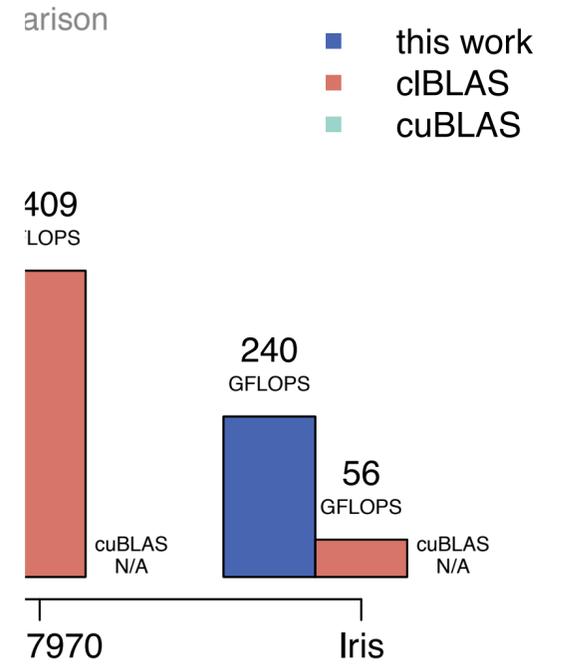
build passing

CLBlast is a modern, lightweight, performant and tunable OpenCL BLAS library written in C++11. It is designed to leverage the full performance potential of a wide variety of OpenCL devices from different vendors, including desktop and laptop GPUs, embedded GPUs, and other accelerators. CLBlast implements BLAS routines: basic linear algebra subprograms operating on vectors and matrices.



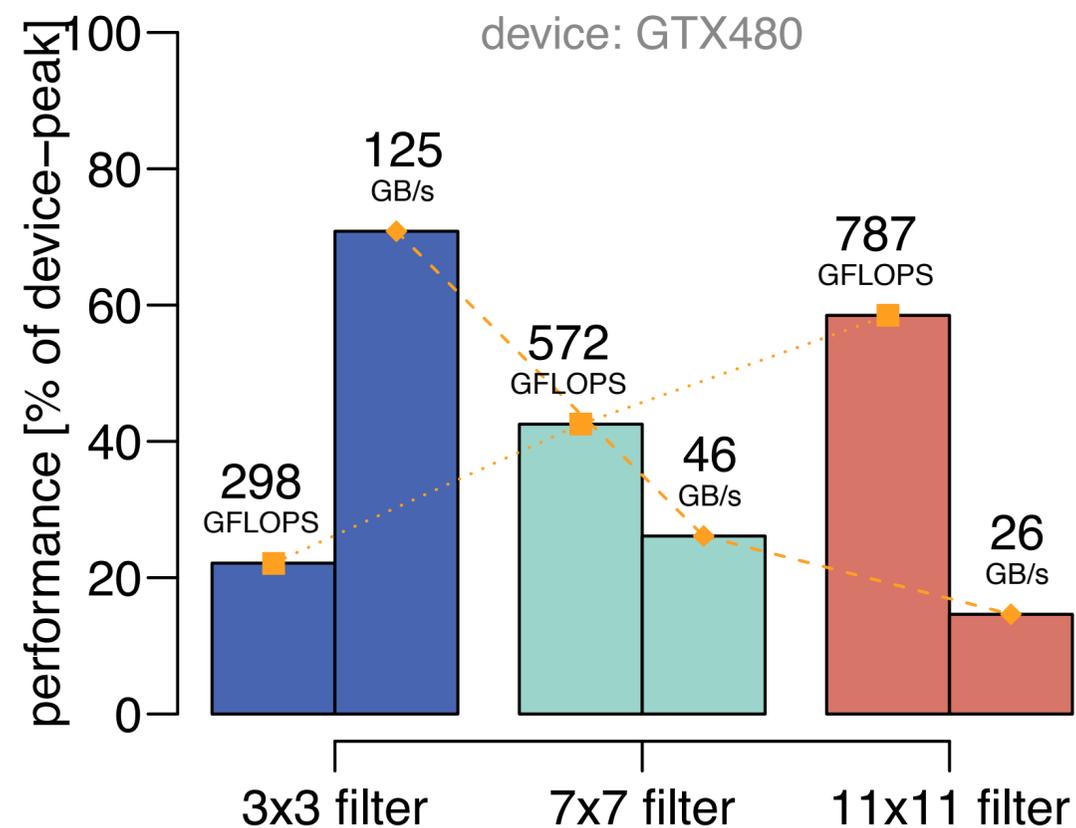
parameter(s)	allowed values
M_{wg}, N_{wg}, K_{wg}	{16,32,64,128}
M_{dimC}, N_{dimC}	{8,16,32}
$L\$A, L\B	{yes,no}
M_{dimA}, N_{dimB}	{8,16,32}
M_{stride}, N_{stride}	{yes,no}
M_{vec}, N_{vec}	{1,2,4,8}
K_{wi}	{2,8}

ameters for
er than cBLAS,
s assembly-



Convolution case-study

applied to a filter of size	best parameters for		
	3x3	7x7	11x11
3x3	100%	82%	64%
7x7	65%	100%	83%
11x11	66%	75%	100%



Conclusions:

- Different best parameters for different:
 - devices (see paper)
 - filter-sizes
- Performance equal or better than the state-of-the-art [3]

[3]: B. Van Werkhoven, J. Maassen, H.E. Bal, and F.J. Seinstra. *Optimizing Convolution Operations on GPUs Using Adaptive Tiling.*

What

[CNugteren / CLCudaAPI](#)
Unwatch 2
Unstar 3
Fork 2

[Code](#)
[Issues 0](#)
[Pull requests 0](#)
[Wiki](#)
[Pulse](#)
[Graphs](#)
[Settings](#)

A portable high-level API with CUDA or OpenCL back-end — Edit

40 commits
 2 branches
 3 releases
 1 contributor

Branch: master
 [New pull request](#)
[New file](#)
[Upload files](#)
[Find file](#)
[SSH](#)
[git@github.com:CNugtere](#)
[Download ZIP](#)

CNugteren Merge pull request #3 from CNugteren/development
 Latest commit 0ad67ce on Nov 1, 2015

cmake/Modules	Changed the name Claduc to CLCudaAPI	6 months ago
doc	Added new method Device::MaxAllocSize to the API	4 months ago
include	Updated to version 4.0	4 months ago
samples	Fixed compiler warnings and errors for Windows using MSVC	4 months ago
test	Fixed bugs for the CUDA unit tests	4 months ago
.gitignore	Initial commit	8 months ago
CHANGELOG	Updated to version 4.0	4 months ago
CMakeLists.txt	Updated to version 4.0	4 months ago
LICENSE	Initial commit	8 months ago
README.md	Fixed compiler warnings and errors for Windows using MSVC	4 months ago

[README.md](#)

CLCudaAPI: A portable high-level API with CUDA or OpenCL back-end

CLCudaAPI provides a C++ interface to the OpenCL API and/or CUDA API. This interface is high-level: all the details of setting up an OpenCL platform and device are handled automatically, as well as for example OpenCL and CUDA memory management. A similar high-level API is also provided by Khronos's `cl.hpp`, so why would someone use CLCudaAPI instead? The main reason is portability: CLCudaAPI provides two header files which both implement the exact same API, but with a different back-end. This allows **porting between OpenCL and CUDA by simply changing the header file!**

CLCudaAPI is written in C++11 and wraps CUDA and OpenCL objects in smart pointers, thus handling memory management automatically. It uses the CUDA driver API, since this is the closest to the OpenCL API, but it uses the OpenCL terminology, since this is the most generic. It compiles OpenCL and/or CUDA kernels at run-time, possible in CUDA only since release 7.0. CLCudaAPI handles the host API only: it still requires two versions of the kernel (although some simple defines could omit this requirement).

ort?

Written for

- Uses the

```
// Copy buf
auto gpuMem
gpuMemory.W
```

- And CLC

- Switch b

- Uses the

vrtc)

Better Than All the Rest:

Finding Maximum-Performance GPU Kernels Using Auto-Tuning

Auto-tuning GPU kernels:

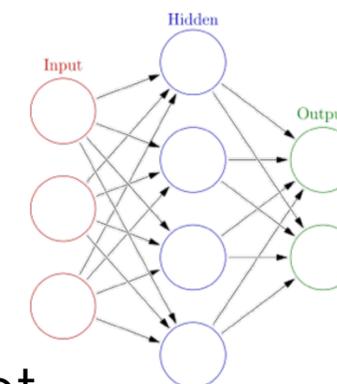
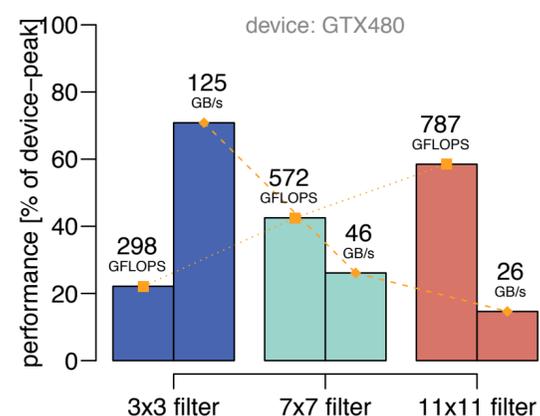
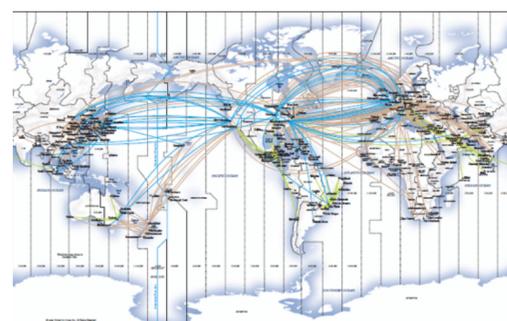
- Large search-space
- Wide variety of devices
- User-parameter dependent

Advanced search strategies:

- Simulated annealing
- Particle swarm optimisation

Case-studies:

- Fastest 2D convolution
- CLBlast matrix-multiplication library



Machine-learning:

- Train a model on a small subset
- Use the model to predict the remainder

Source-code on GitHub:
<https://github.com/CNugteren/CLTune>

[More info]: C. Nugteren and V. Codreanu. CLTune: A Generic Auto-Tuner for OpenCL Kernels. IEEE MCSoc '15.

Slides available @ www.cedricnugteren.nl