



HPC center of the Netherlands

# CLTune: A Generic Auto-Tuner for OpenCL Kernels

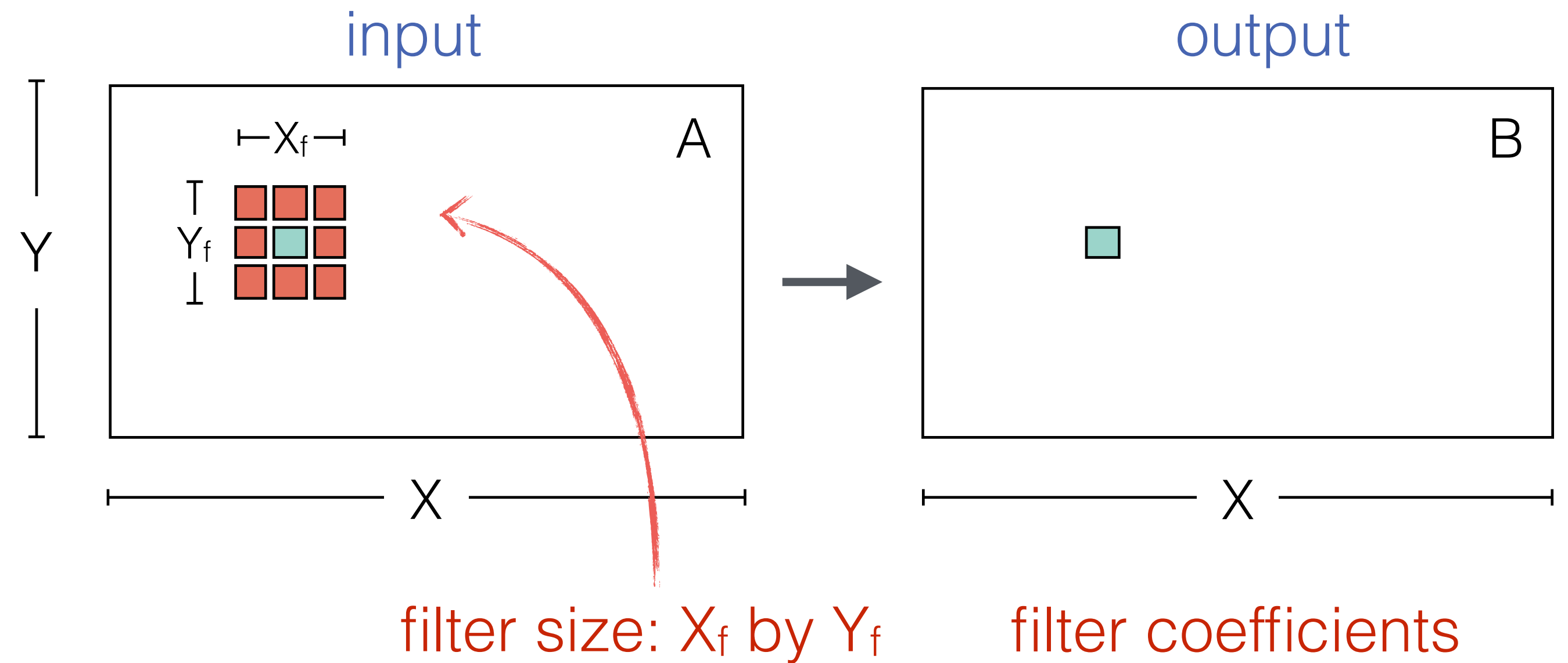
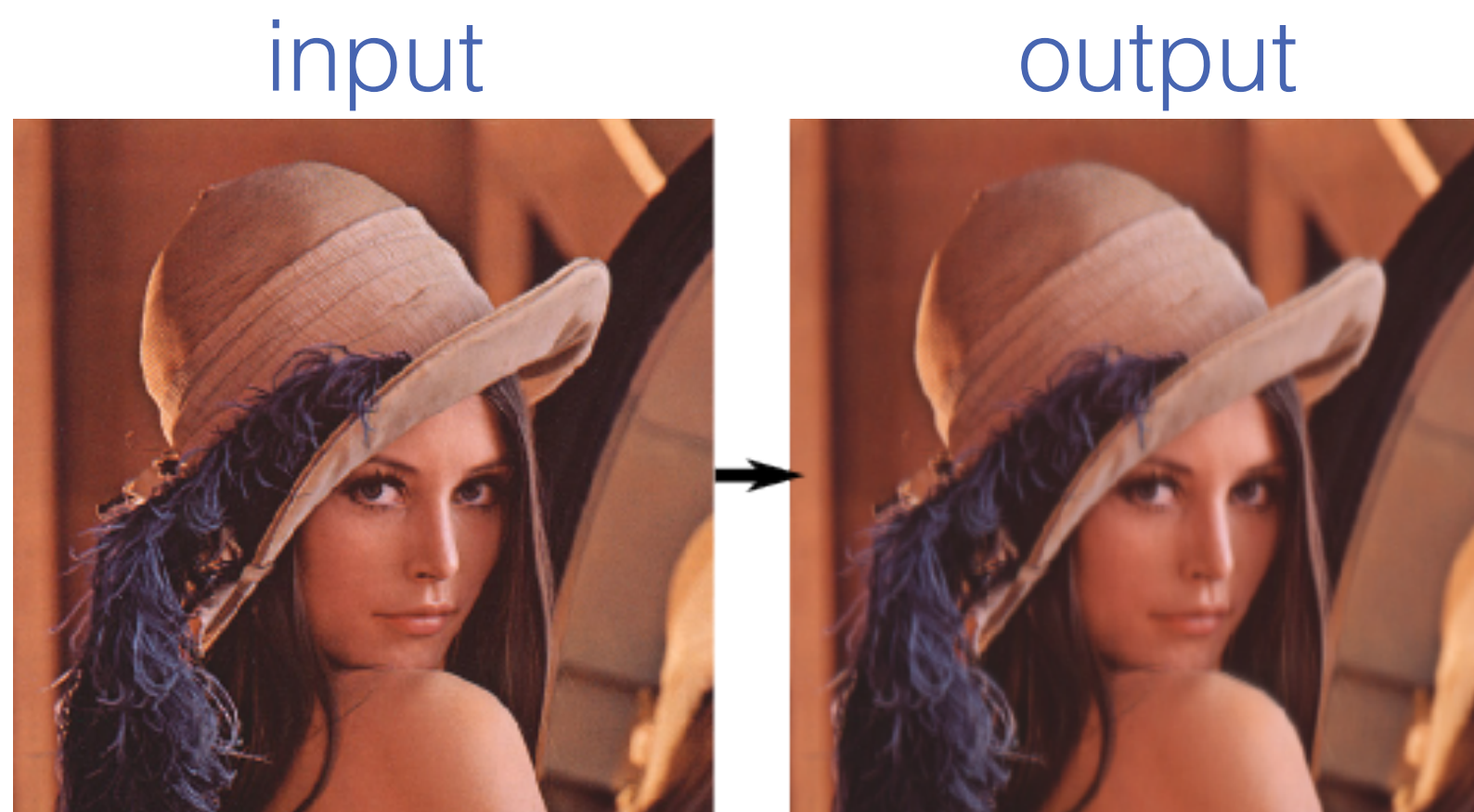
**IEEE MCSoc**

September 24, 2015

Cedric Nugteren (presenter), Valeriu Codreanu

# Example: convolution

## Example: blur filter



## Targets:

- GPUs
- Multi-core CPUs
- Other OpenCL-capable devices

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

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 }
```

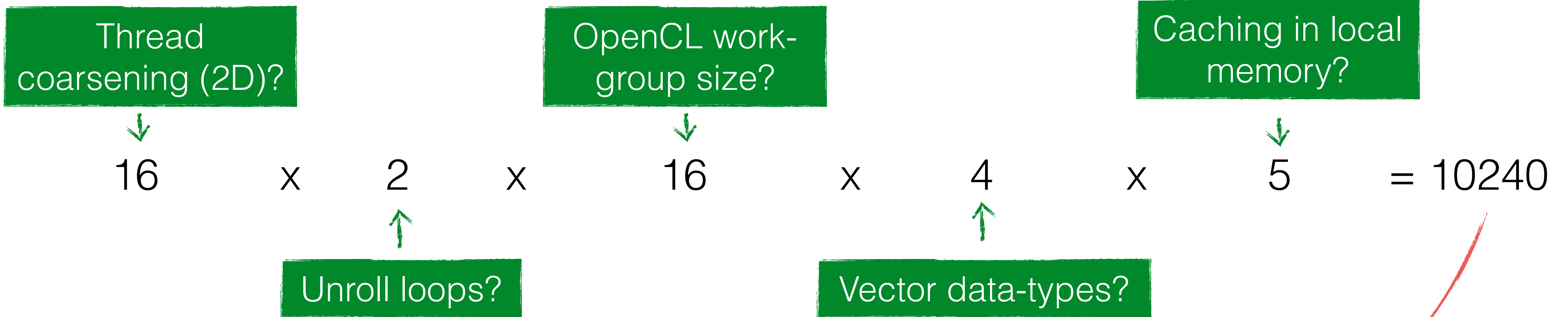
OpenCL work-group size?

Vector data-types?

Caching in local memory?



# Search-space explosion



## Large search-space:

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

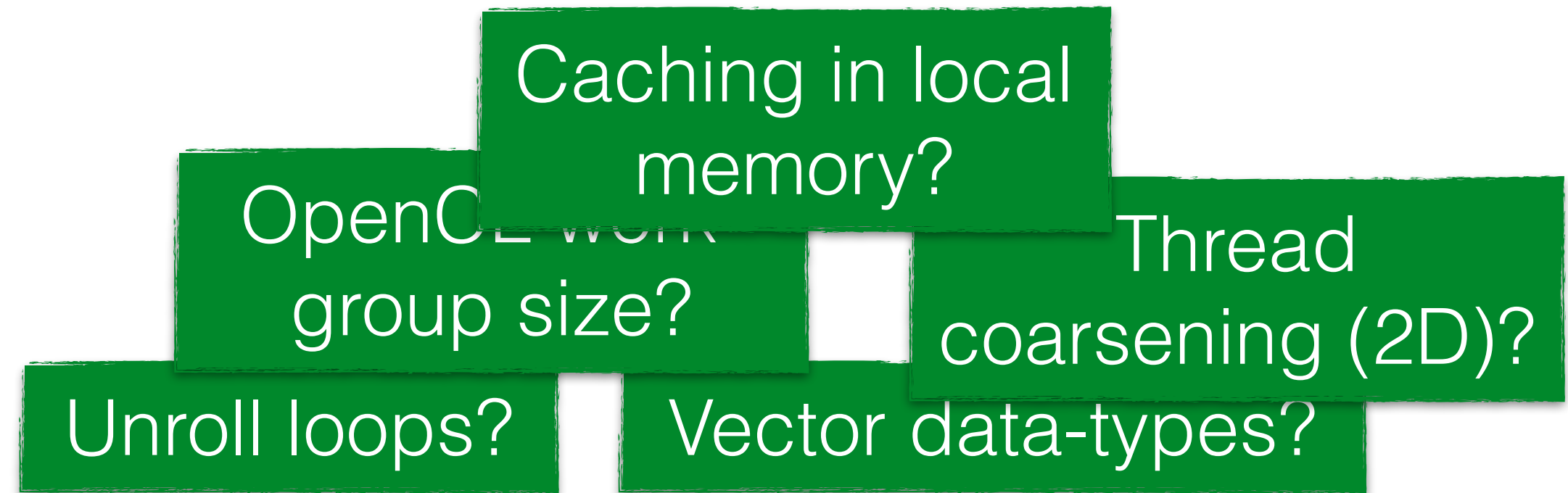
3424 configurations

filter illegal configurations

# Why do we need an auto-tuner?

## Large search-space:

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



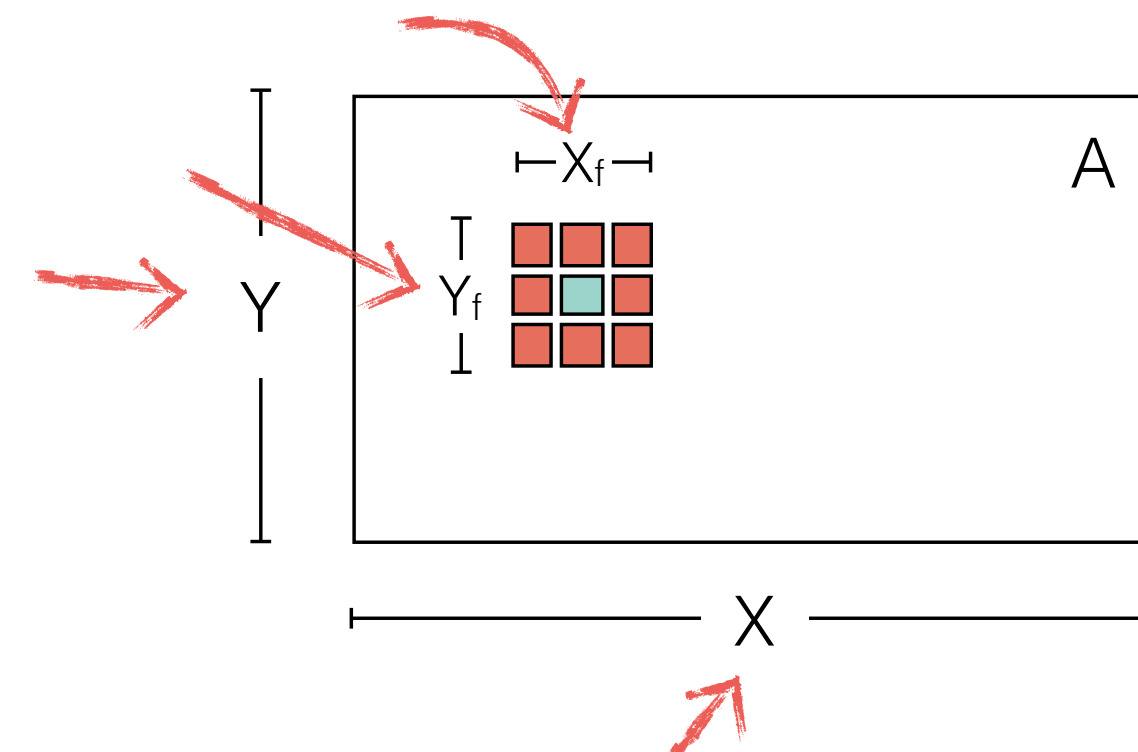
## Wide variety of devices:

- Different optimal kernels
- Even from the same vendor

vendor and device name	architecture	compiler and SDK	peak GFLOPS	peak GB/s	GFLOPS per GB/s
NVIDIA Tesla K40m	Kepler	CUDA 7.0	4291	288	14.9
NVIDIA GeForce GTX480	Fermi	CUDA 5.5	1345	177	7.6
AMD Radeon HD7970	Tahiti	APP 2.9	4368	288	15.1
Intel Iris 5100	Iris	Apple 2.4.2	832	26	32.5

## User-parameter dependent:

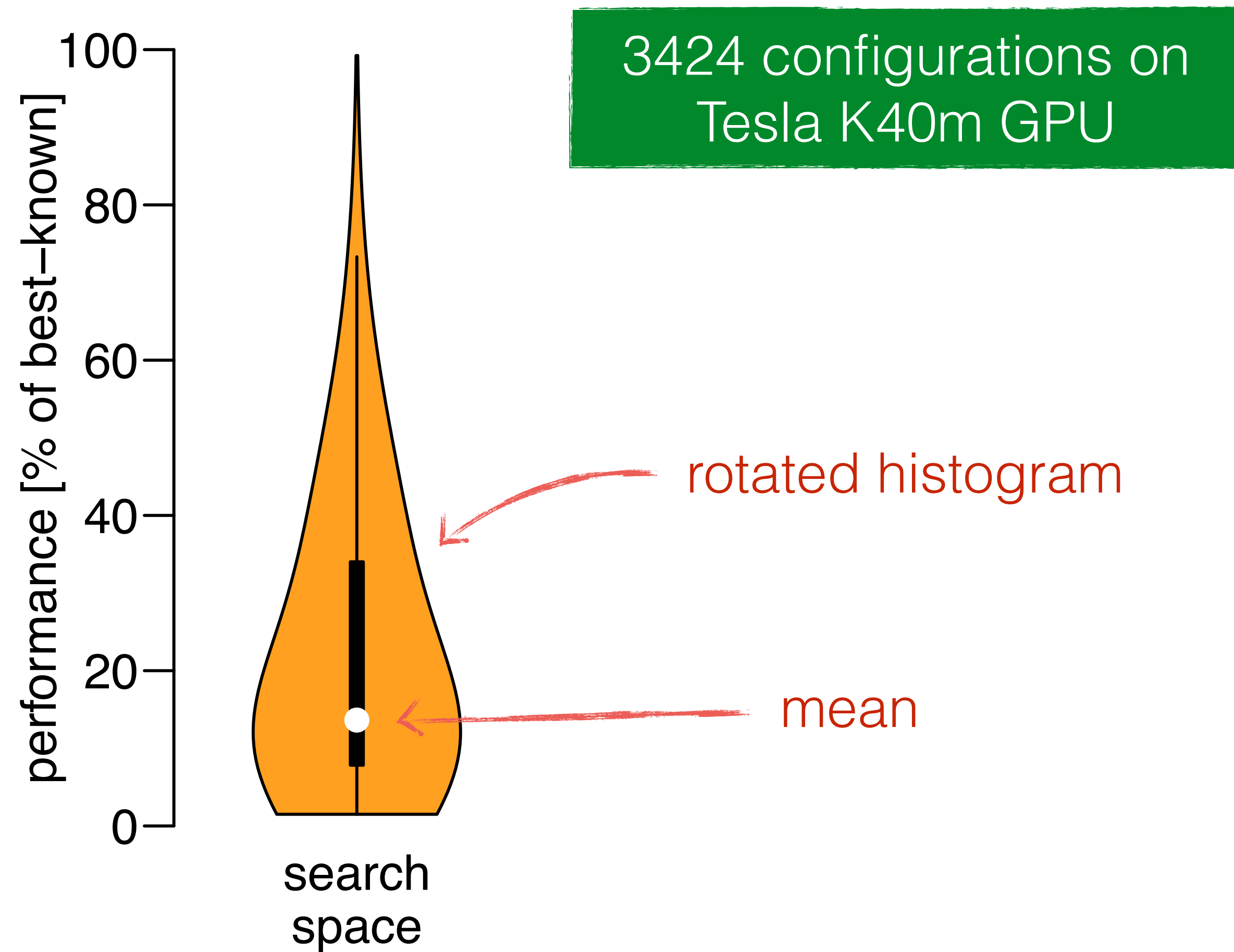
- Examples: matrix sizes, image size, filter sizes, etc.



# Search strategies

## Option 0: Full search

- 😊 Finds optimal solution
- 😞 Explores all options



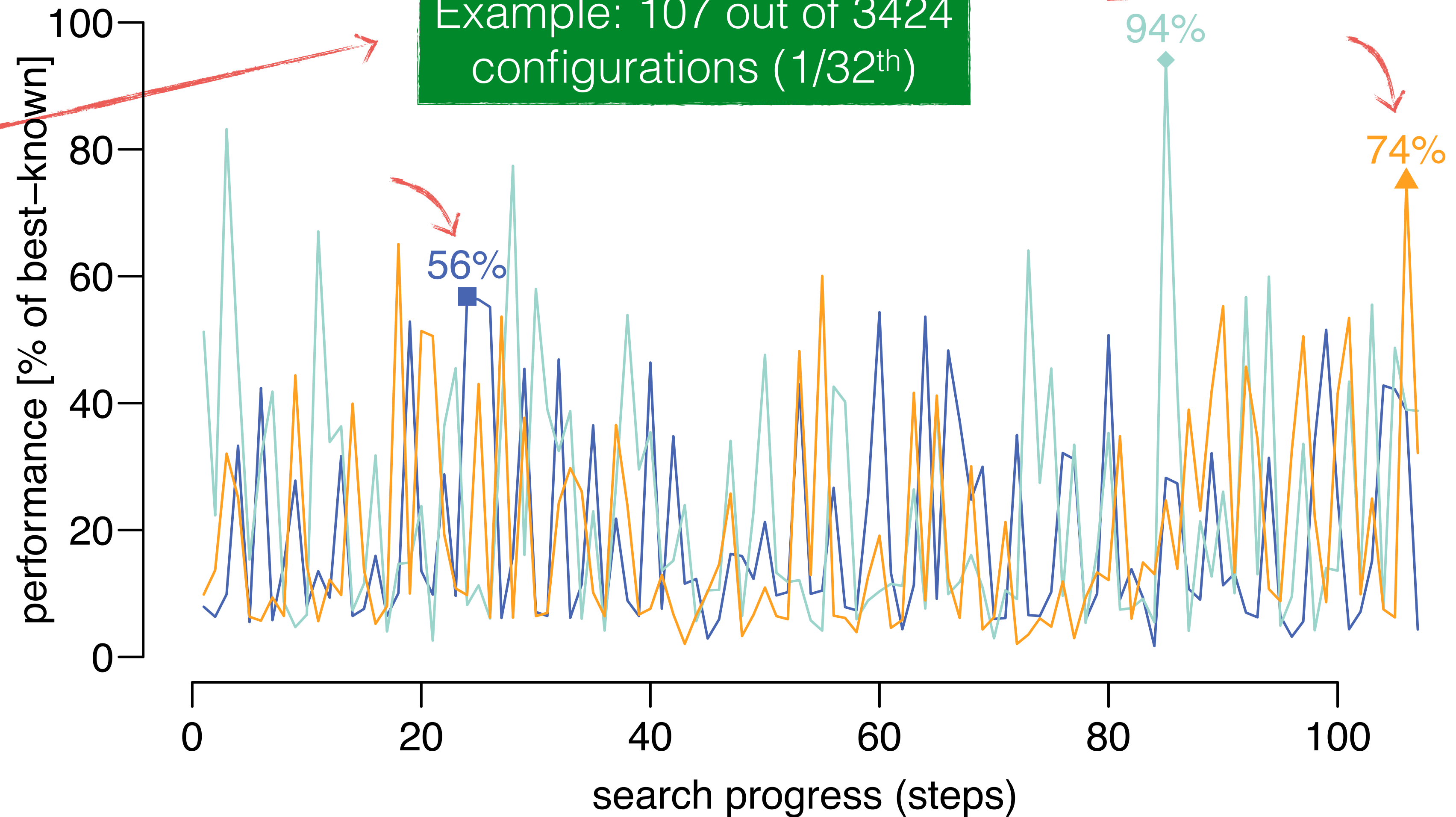
# Search strategies

**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**

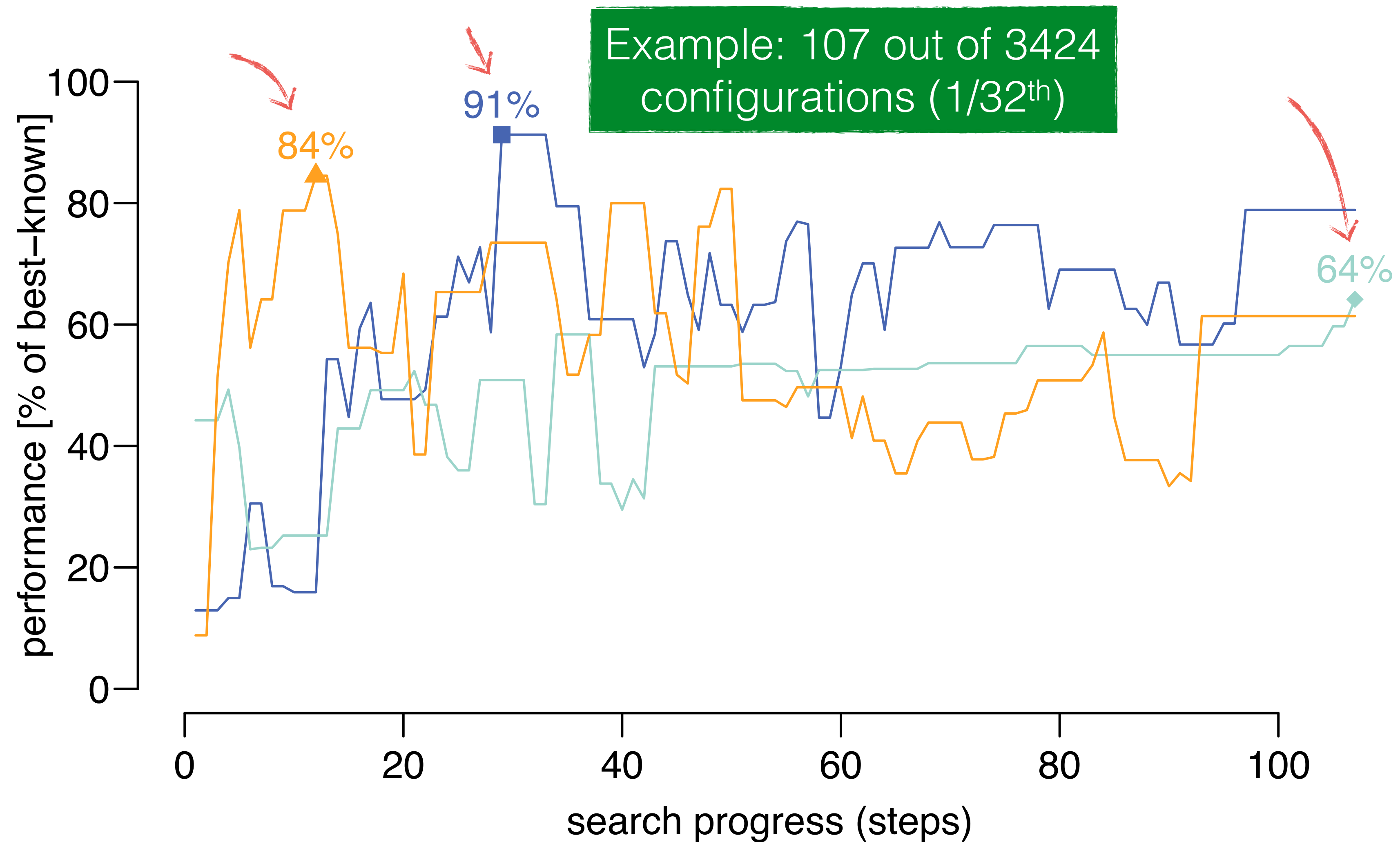
😊 Explores arbitrary fraction

☹ 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 optimisation**

😊 Explores arbitrary fraction

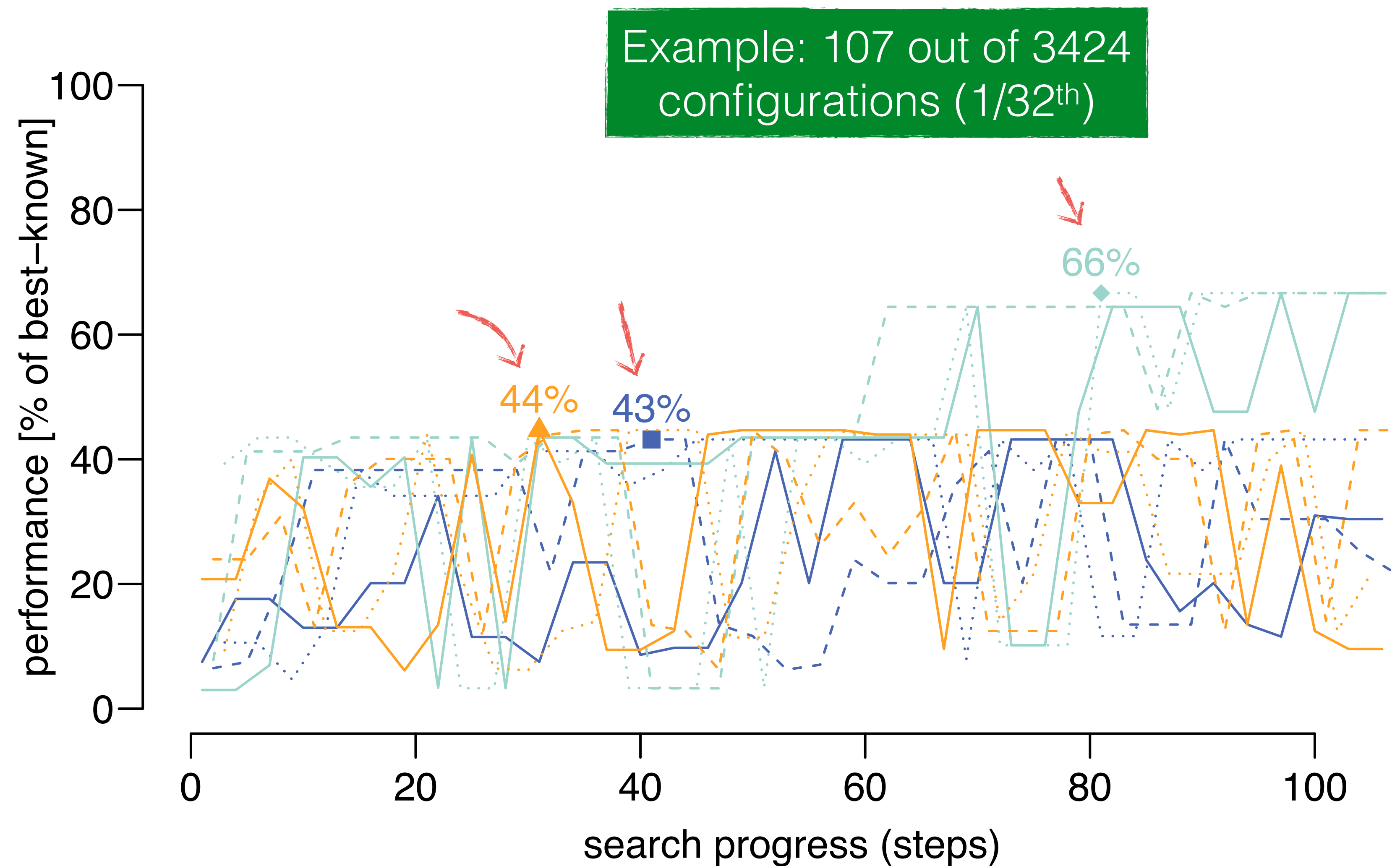
☹ Performance varies

☹ Meta-parameter

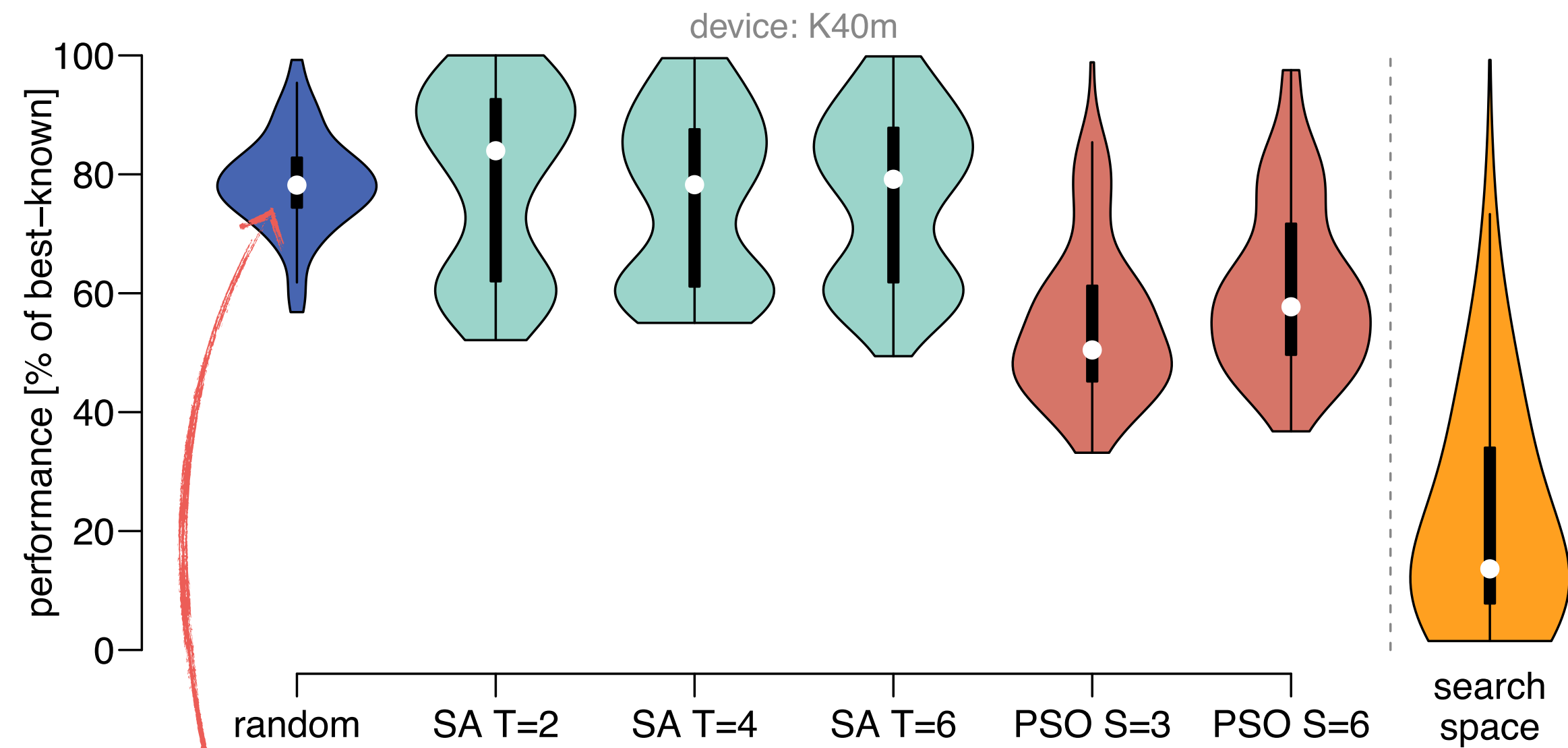
☹ Local optima

Colours: 3 example runs

Line-types: 3 swarms



# Search strategies evaluation

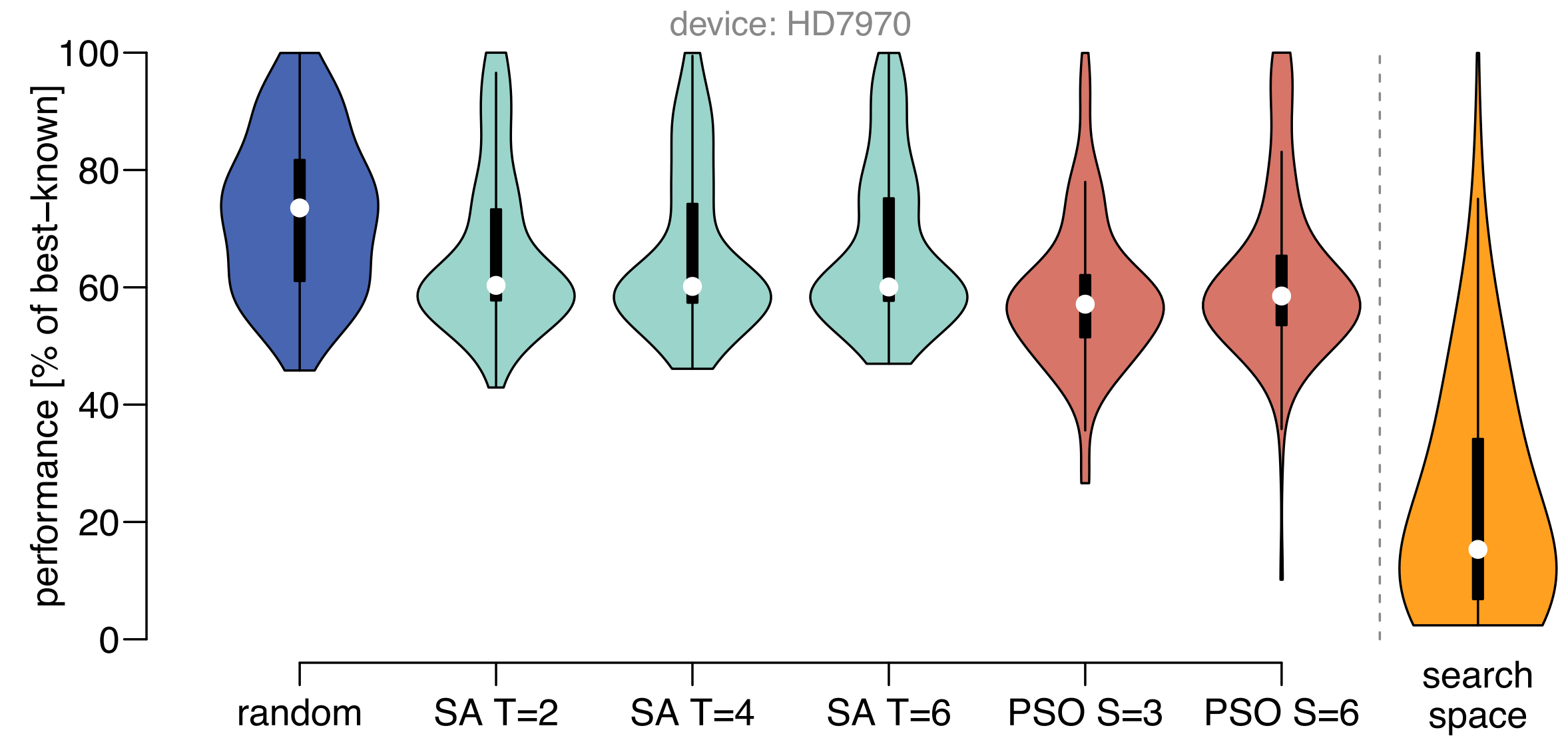
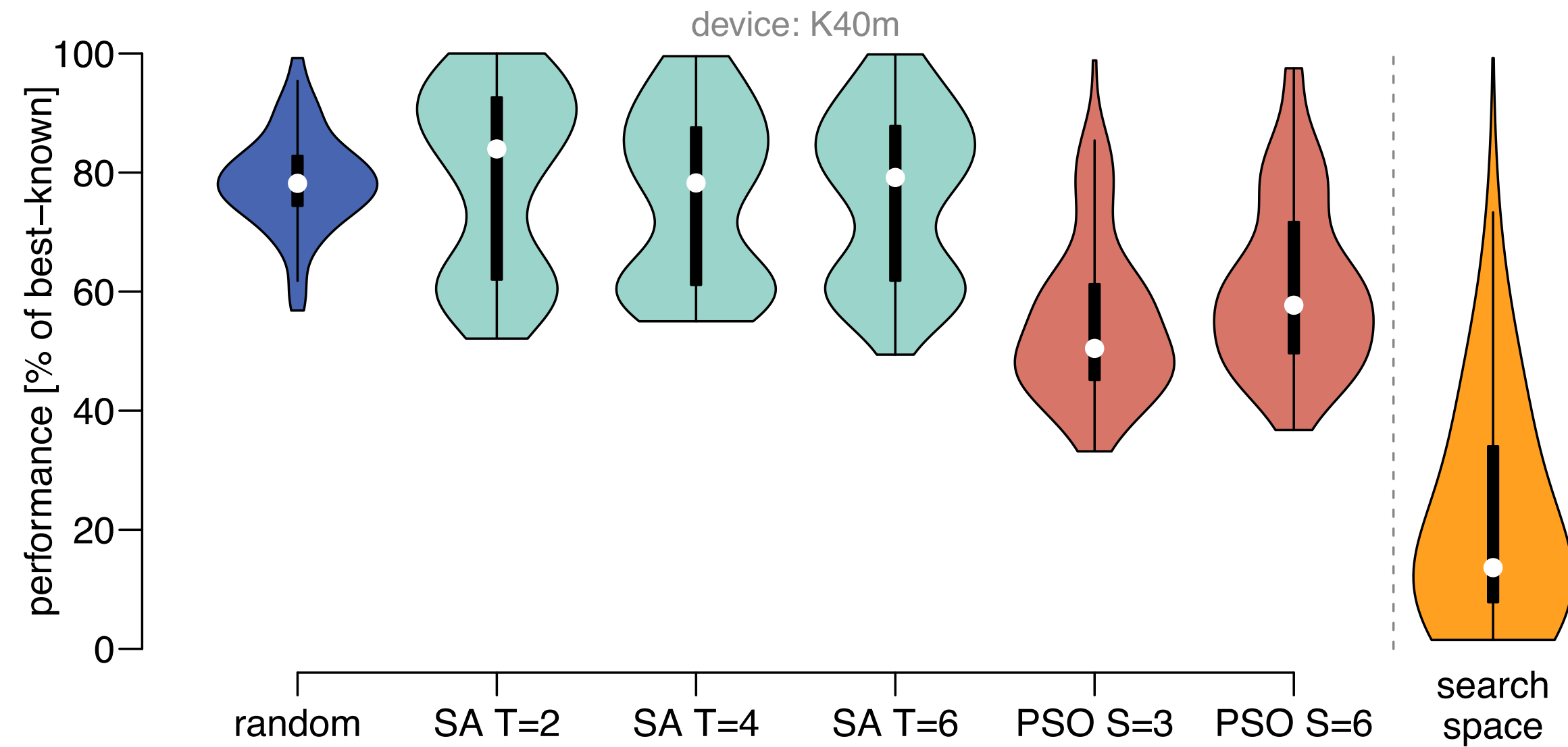


average best result  
of 128 searches

meta-parameters  
for SA and PSO

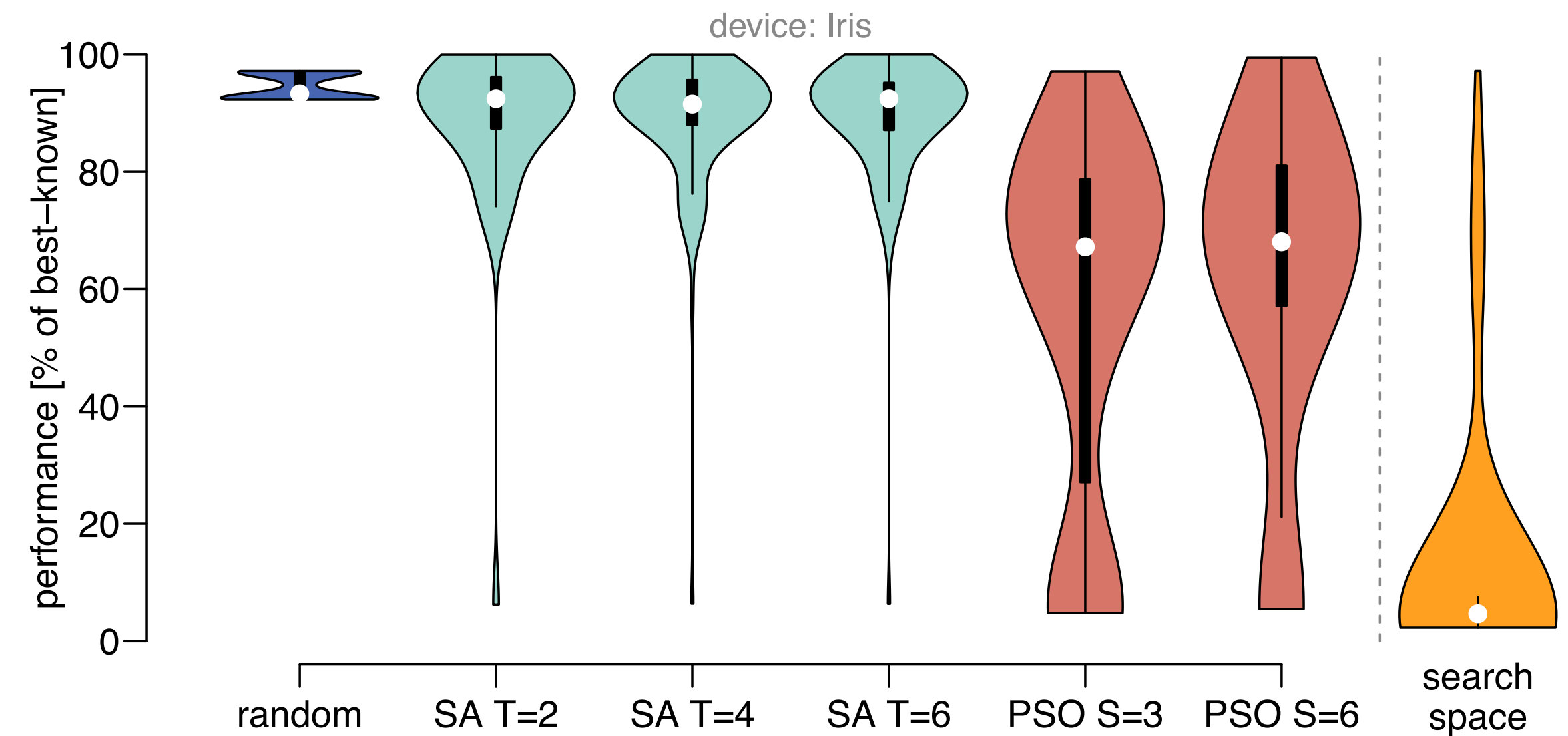
Each search: 107 out of 3424  
configurations ( $1/32^{\text{th}}$ )

# Search strategies evaluation



## Conclusions:

- Different per device
- PSO performs poorly
- Random search and SA perform well



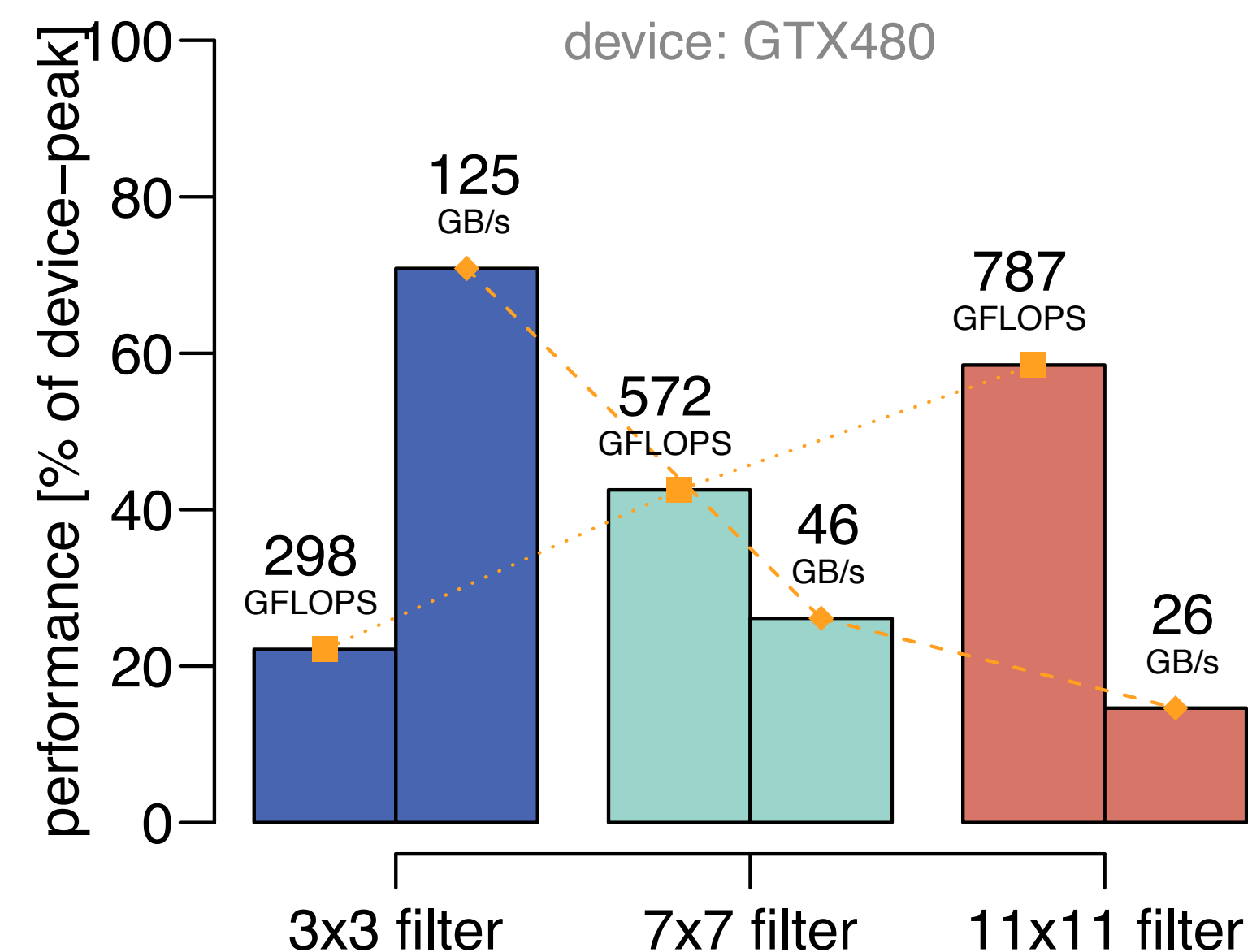
# Convolution case-study

parameter(s)	allowed values	GeForce GTX480		
		3x3	7x7	11x11
$X_{wg}, Y_{wg}$	{8,16,32,64}	64,8	32,8	32,8
$X_{wpt}, Y_{wpt}$	{1,2,4,8}	1,4	2,8	2,4
$L\$$	{0,1,2}	0	2	1
$VW$	{1,2,4,8}	1	2	2
$PAD$	{0,1}	0	0	0
$UNR$	{yes,no}	yes	yes	yes

## Conclusions:

- Different best parameters for different:
  - devices (see paper)
  - filter-sizes
- Performance equal or better than the state-of-the-art [1]

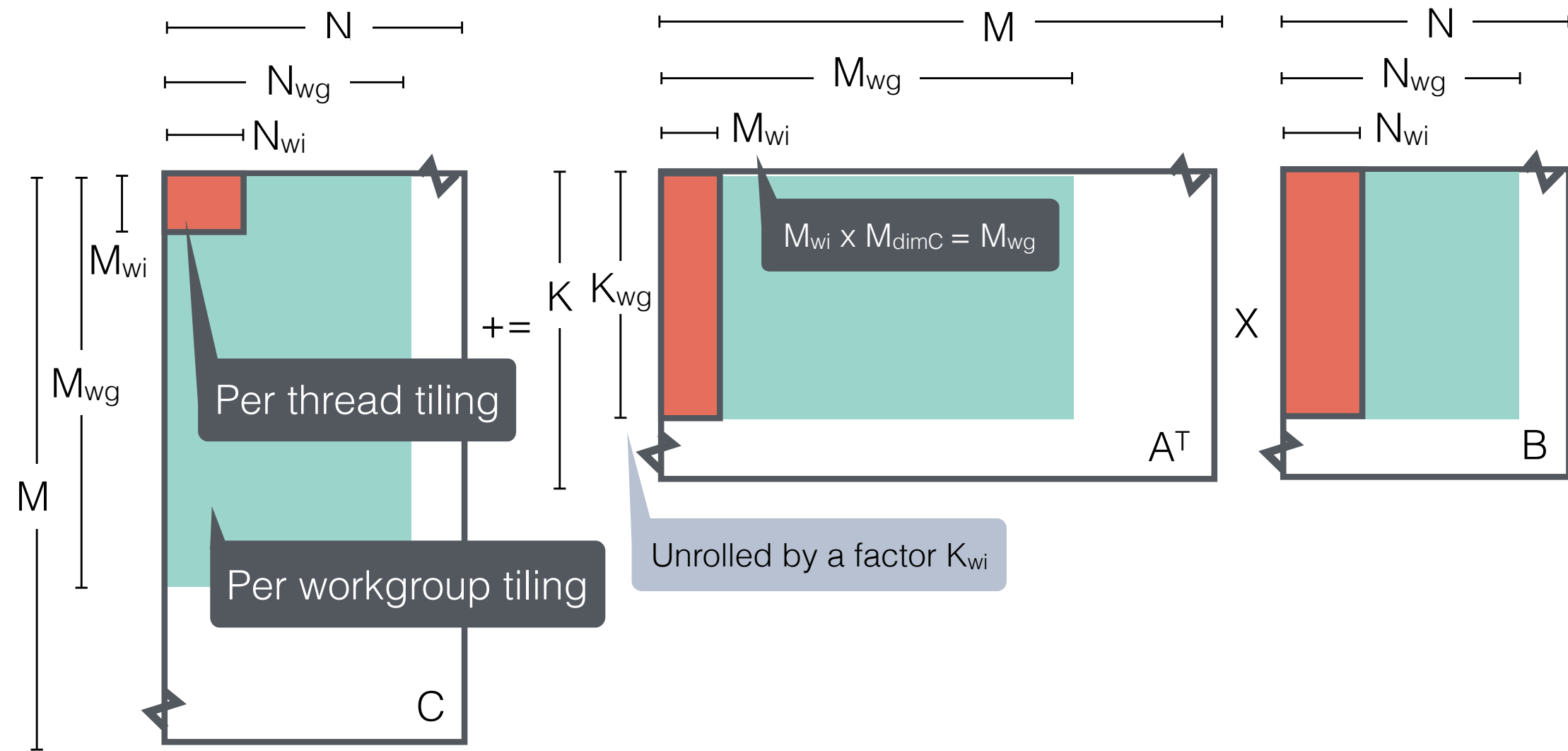
applied to a filter of size	best parameters for		
	3x3	7x7	11x11
3x3	100%	82%	64%
7x7	65%	100%	83%
11x11	66%	75%	100%



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



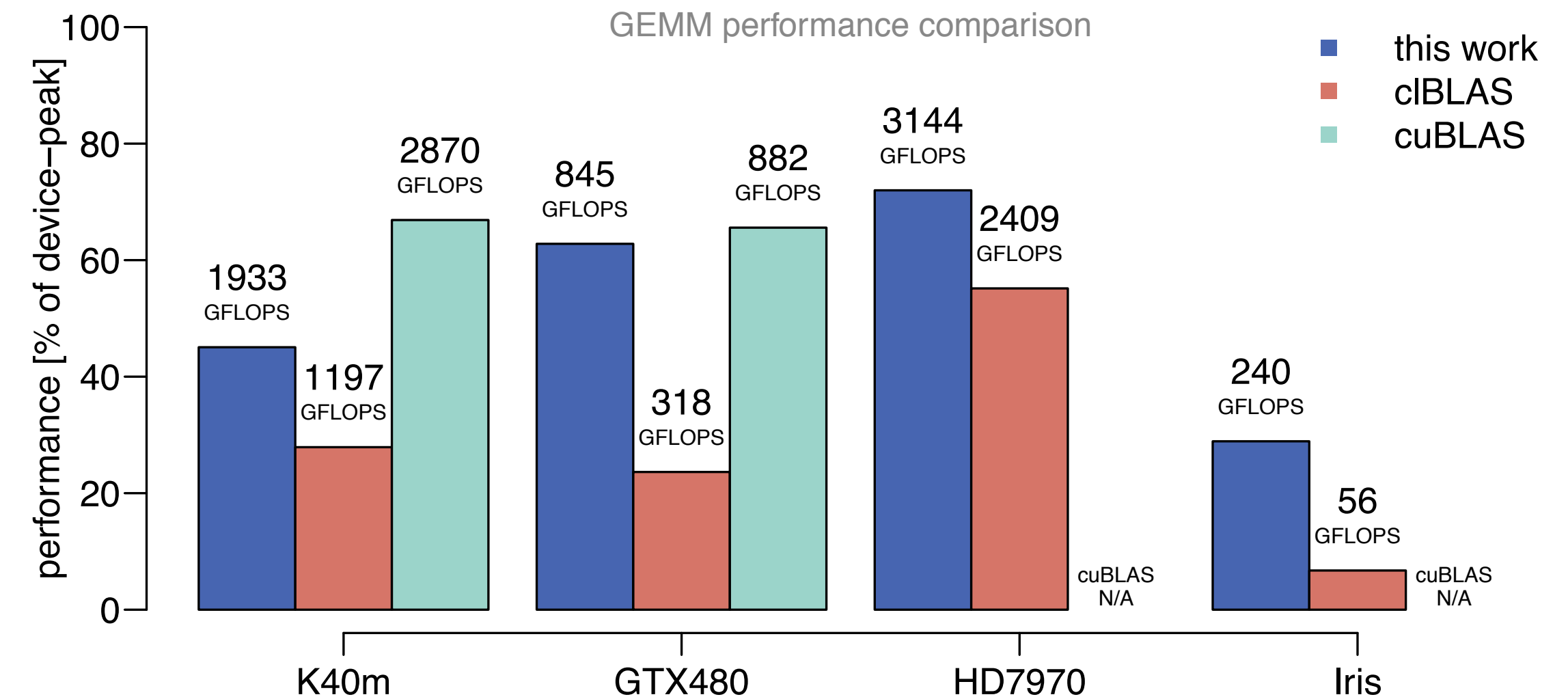
# GEMM case-study



## 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



# CLTune: A Generic Auto-Tuner for OpenCL Kernels

## Auto-tuning OpenCL kernels:

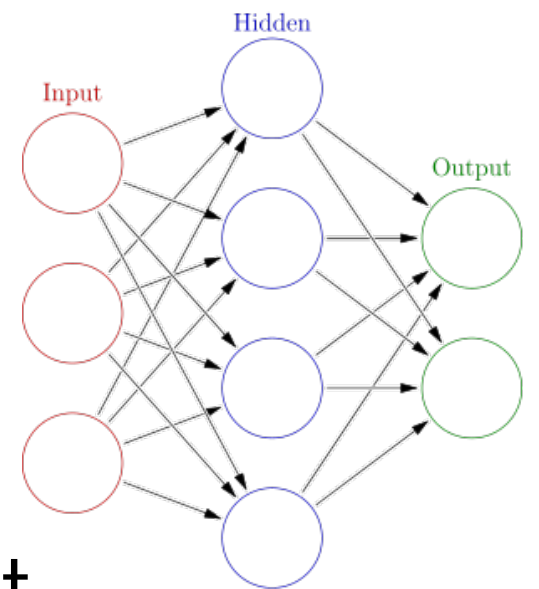
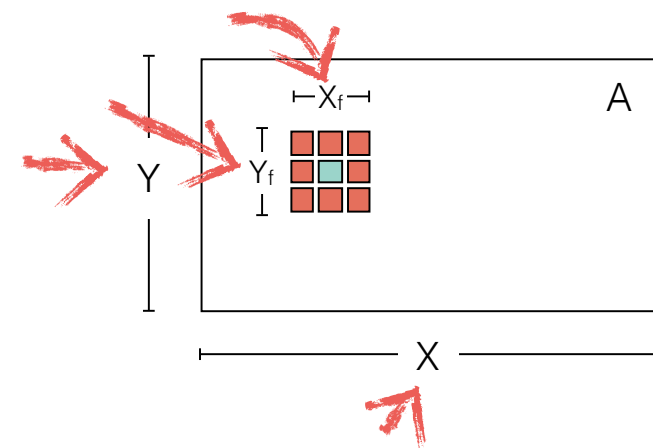
- Large search-space
- Wide variety of devices
- User-parameter dependent

## Advanced search strategies:

- Simulated annealing
- Particle swarm optimisation

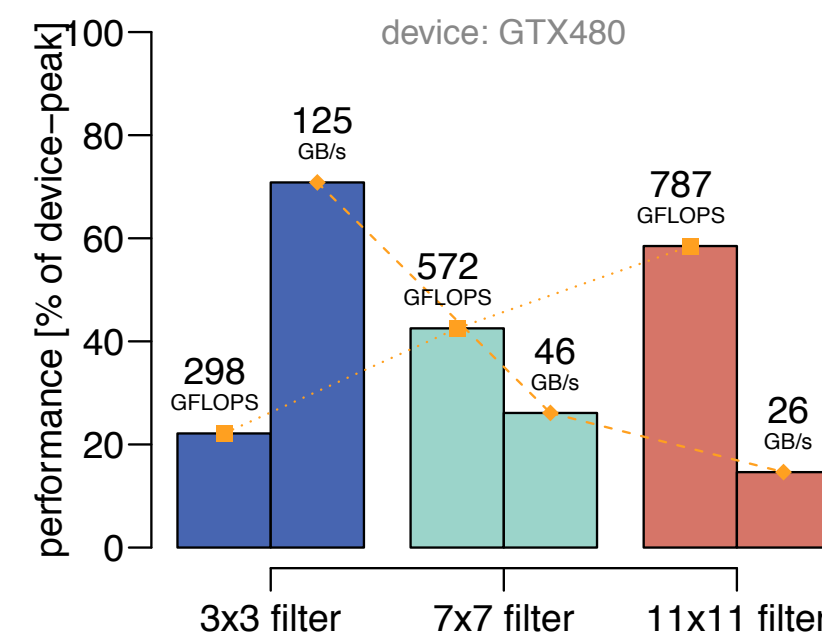
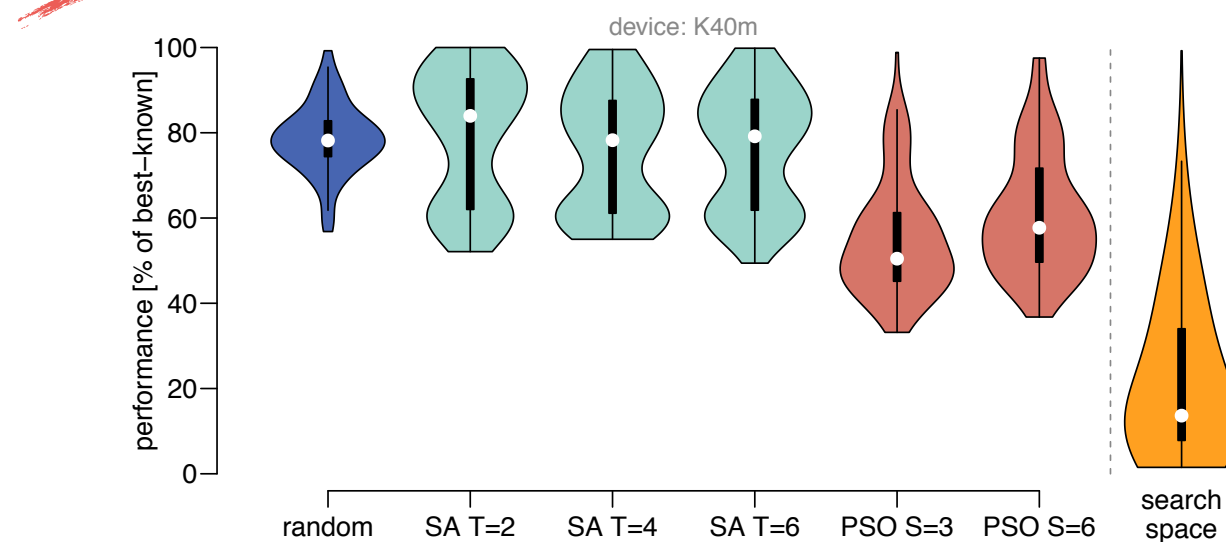
## Case-studies:

- Fastest 2D convolution
- Fast matrix-multiplication



## Future: machine-learning [2]

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



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

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