



Efficient Auto-Tuning of Parallel Programs with Interdependent Tuning Parameters via Auto-Tuning Framework (ATF)

Ari Rasch, Richard Schulze, Michel Steuer,
Sergei Gorlatch



THE UNIVERSITY
of EDINBURGH

Please Note

This talk will be on a quite high level:

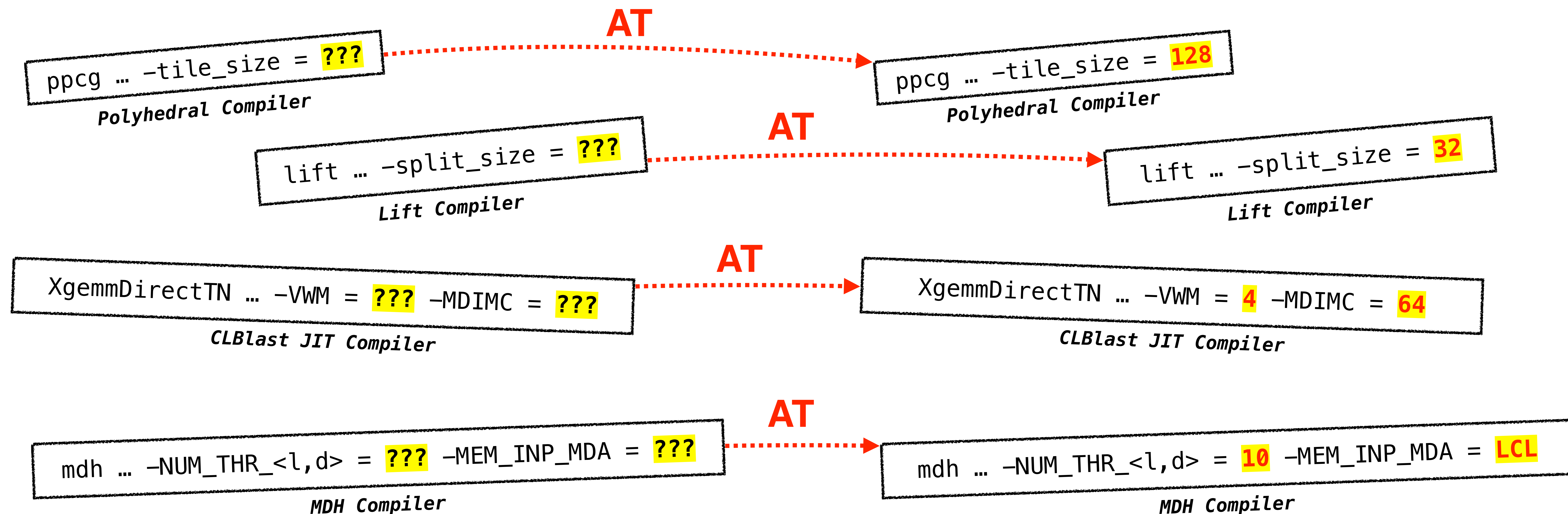
- **focus of this talk are the “what” & “why” questions;**
- **we address the “how” question by illustrating our basic ideas only;**
- **details about our approach can be found in the paper.**

What is *Auto-Tuning*?

Auto-Tuning (AT) aims at automatizing the process of code optimization:

Auto-Tunable Optimization Process

Auto-Tuned Optimization Process

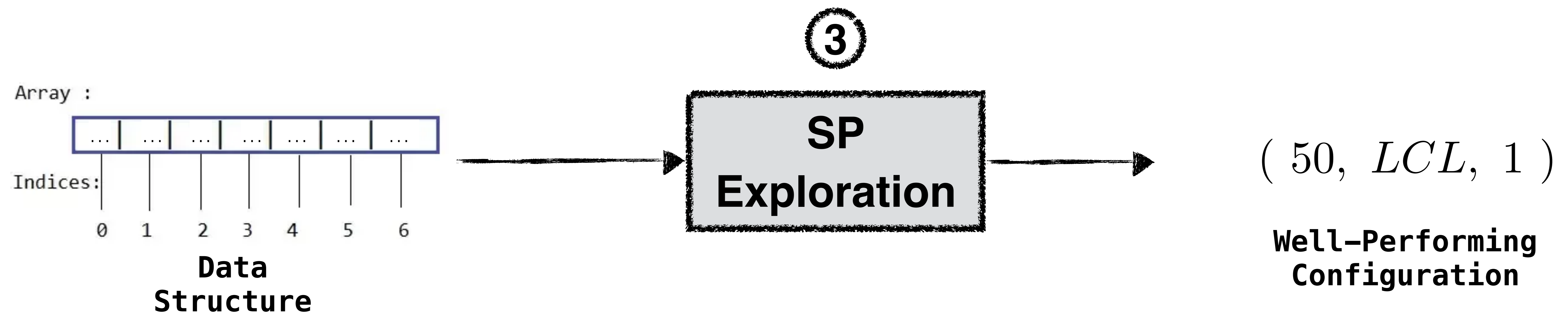
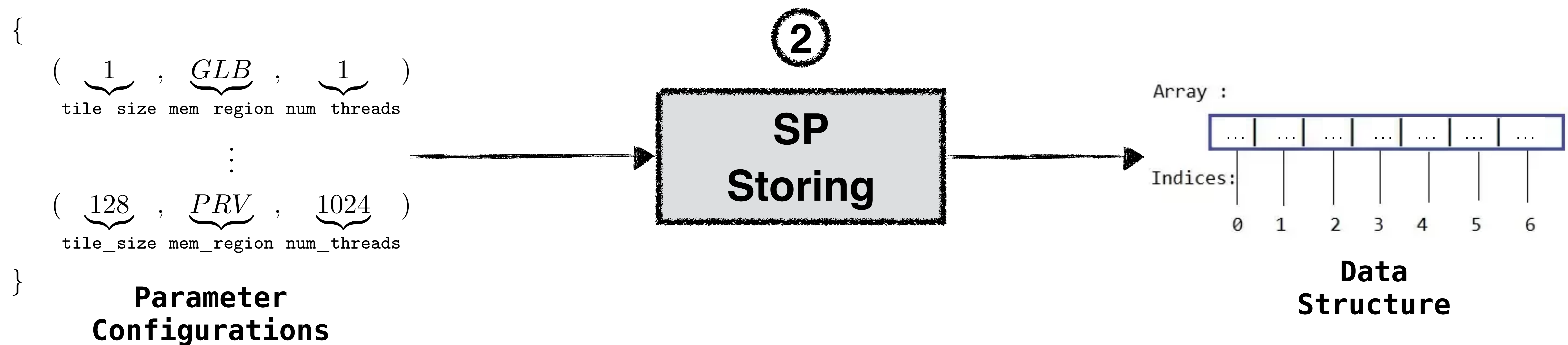
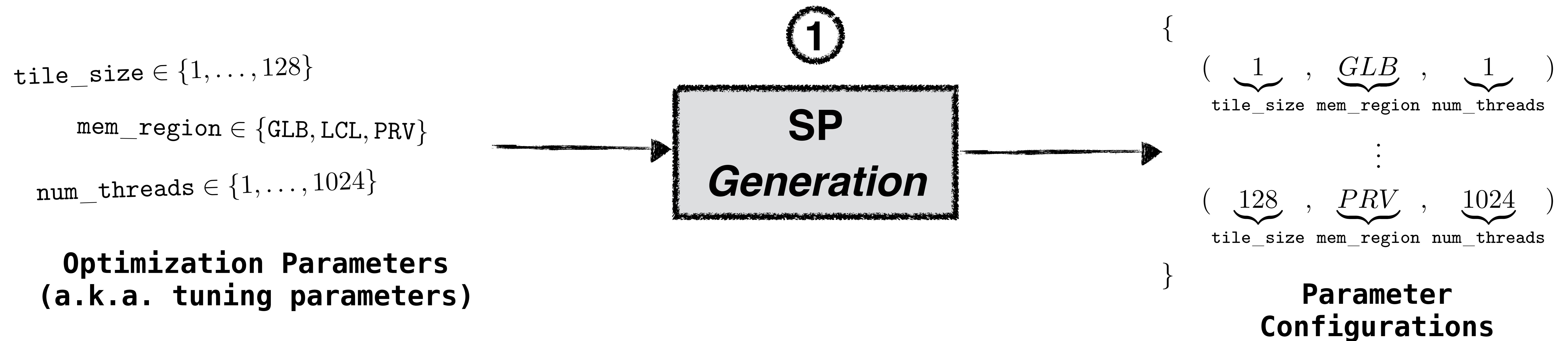


Generic Optimization Parameters

Instantiated Optimization Parameters

What is *Auto-Tuning*?

Auto-Tuning usually consists of three major phases:



What is *Auto-Tuning*?

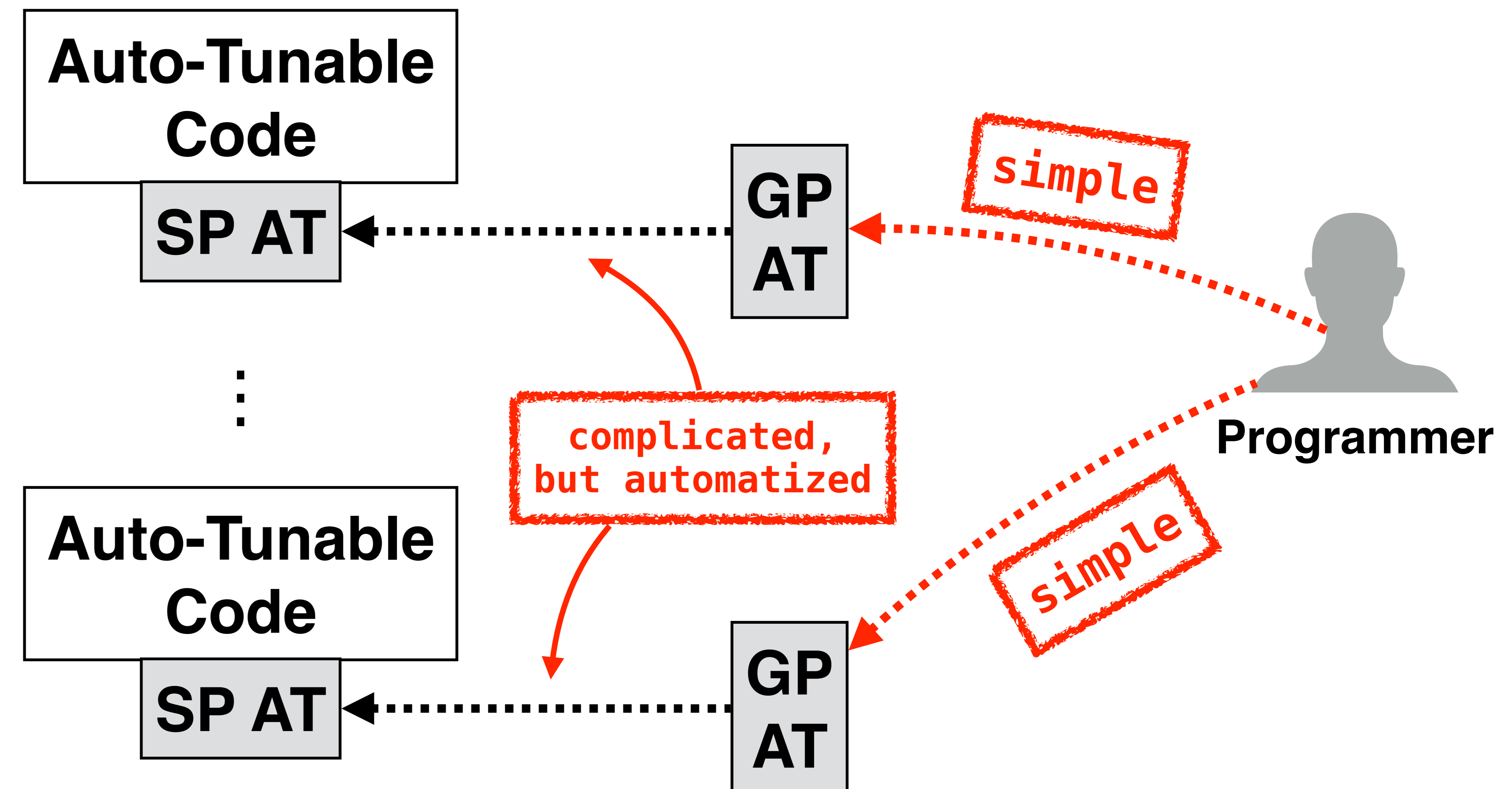
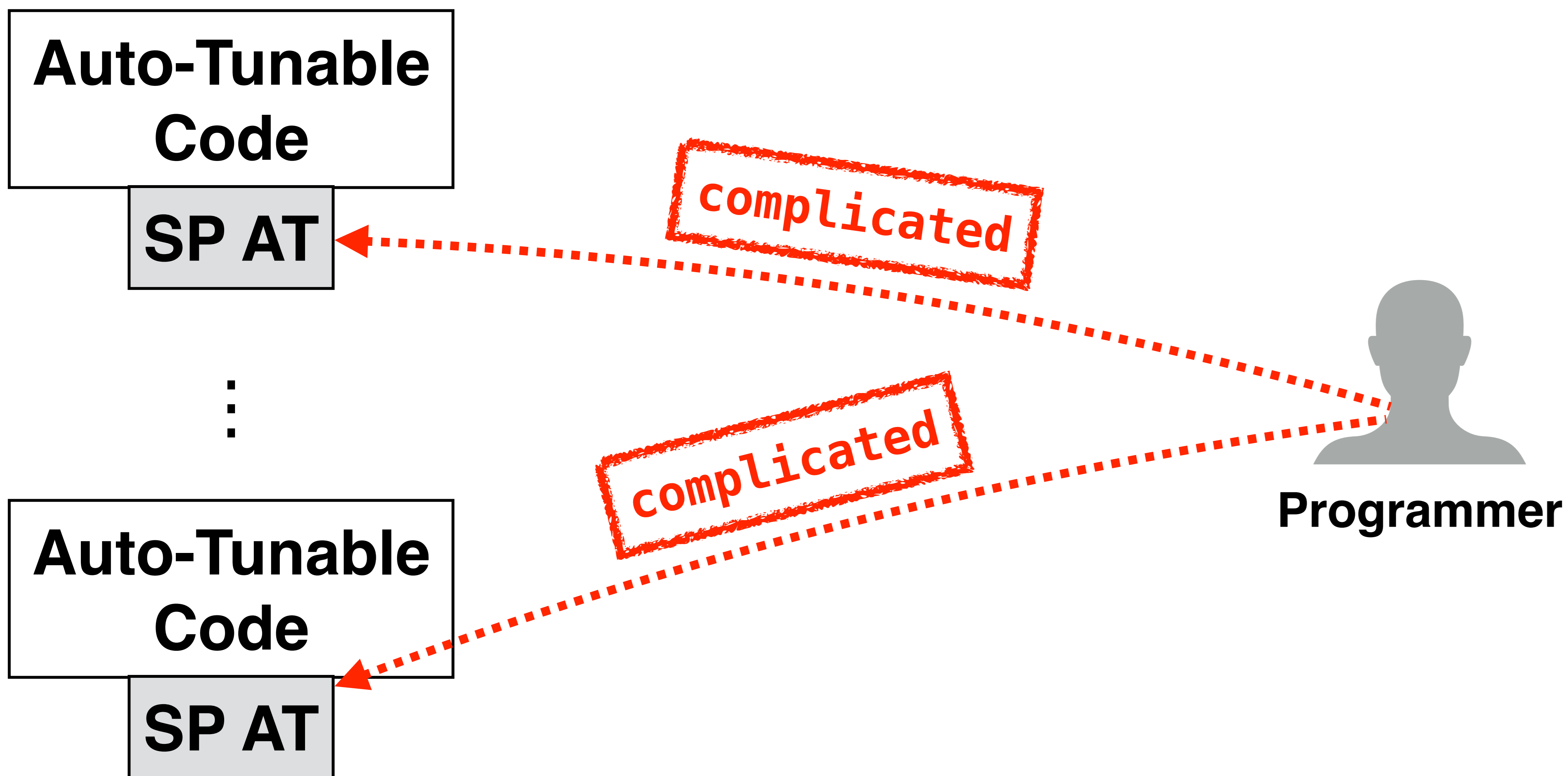
Auto-Tuning (AT) can be categorized into two major categories:

Special-Purpose (SP) Auto-Tuning

General-Purpose (GP) Auto-Tuning

Basic Idea:

Basic Idea:



Notable Examples:

Notable Examples:

ATLAS FFTW MilePost CHILL
 PATUS SPIRAL Nitro PolyMage Apollo

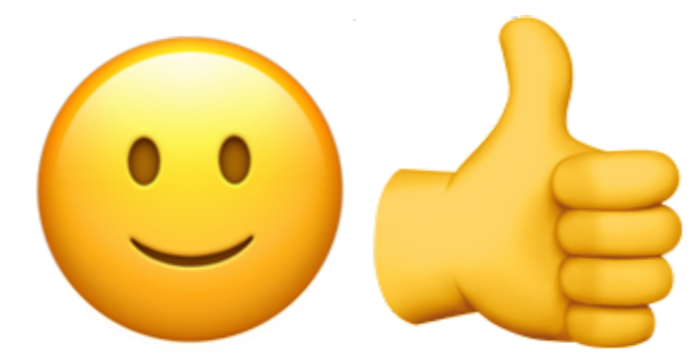
OpenTuner libTuning Orio
 KernelTuner ActiveHarmony
 CLTune

What is *Auto-Tuning*?

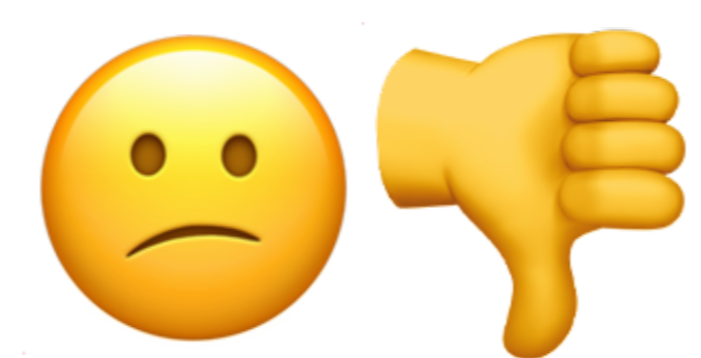
Auto-Tuning (AT) can be categorized into two major categories:

Special-Purpose Auto-Tuning

Special-Purpose Auto-Tuners usually achieve good tuning results

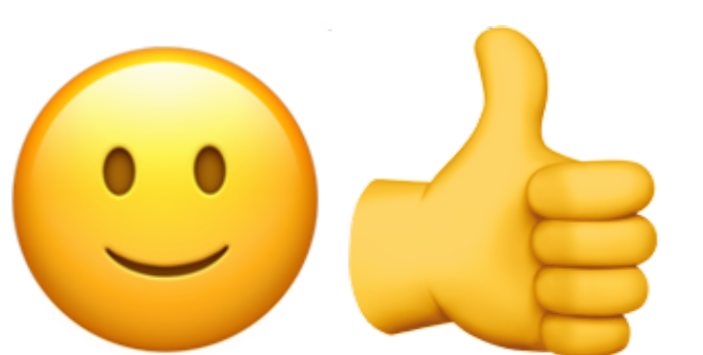


However: they have to be designed & implemented from scratch for each new target application, which requires expert knowledge and is cumbersome

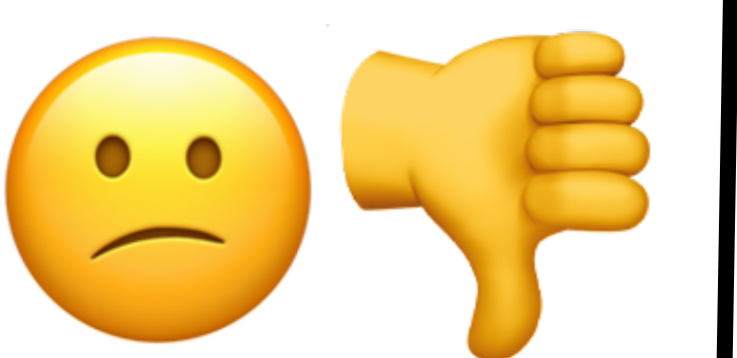


General-Purpose Auto-Tuning

General-Purpose Auto-Tuners automatically generate special-purpose tuners



However: current approaches struggle with programs that have *interdependent* tuning parameters



What are *Interdependent Tuning-Parameters*?

Independent Parameters

Parameters:

$\text{tile_size} \in \{1, \dots, 128\}$

$\text{mem_region} \in \{\text{GLB}, \text{LCL}, \text{PRV}\}$

$\text{num_threads} \in \{1, \dots, 1024\}$

Constraints:

<none>

Configurations:

$\{ (1, \text{GLB}, 1), \dots, (128, \text{PRV}, 1024) \}$

Each combination of parameters' values represents a valid parameter configuration

Interdependent Parameters

Parameters:

$\text{tile_size_1} \in \{1, \dots, 128\}$

$\text{tile_size_2} \in \{1, \dots, 128\}$

$\text{tile_size_3} \in \{1, \dots, 128\}$

Constraints:

$\text{tile_size_3} \mid \text{tile_size_2} \mid \text{tile_size_1}$

Configurations:

$\{ (1, 1, 1), (1, \text{X}, 2), \dots, (128, 128, 128) \}$

*tile_size_2
not multiple of
tile_size_3*

Only combinations that satisfy the constraints represent valid configurations

Interdependent Tuning-Parameters in General-Purpose Auto-Tuning

Current approach have either **no support** or only **limited efficiency** for programs with interdependent tuning parameters:

OpenTuner

Orio

libTuning

no support: invalid configurations are kept in the search space, which severely hinders the tuners from finding well-performing configurations

ActiveHarmony

CLTune

KernelTuner

limited efficiency: the approaches are efficient for small search spaces only, because of sub-optimal process to generating, storing, and exploring the search spaces of interdependent parameters

Goal of this Work

We present three new contributions of our general-purpose Auto-Tuning Framework (ATF):

1. ATF *generates* the search space of interdependent parameters with high performance

2. ATF *stores* such spaces with low memory footprint

3. ATF *explores* these spaces efficiently

Parameter Constraints

Chain-of-Trees Structure

1. Generation

ATF relies on *parameter constraints*, rather than traditional *search space constraints*:

```

for ( v1 : r1 )
  ∴
  for ( vk : rk )
    if( sc(v1, ..., vk) )
      add_config( v1, ..., vk );
  
```

Search Space Constraint (SC)

PCs enable generating groups of interdependent parameters independently & in parallel

```

parallel_for ( G : {G1, ..., Gn} )
{
  parallel_for ( v1G : r1G )
  {
    if( pc1G(v1G) )
    {
      ∴
      parallel_for ( vtgG : rtgG )
      {
        if( pctgG(vtgG) )
        {
          for ( vtg+1G : rtg+1G )
          {
            if( pc(vtg+1G) )
            {
              ∴
              for( vkG : rkG )
              {
                if( pc(vkG) )
                {
                  add_config( v1G, ..., vkG );
                }
              }
            }
          }
        }
      }
    }
  }
}
  
```

Parameter Constraint (PC)

Classical Approach

PCs enable generating individual groups in parallel

PCs enable checking constraints early in the loop nest

- ATF's parameter constraints enable:**
1. Checking Constraints Early
 2. Generating groups of interdependent parameters independently & in parallel
 3. Generating individual groups in parallel

ATF's Approach

2. Storing

Basic Idea

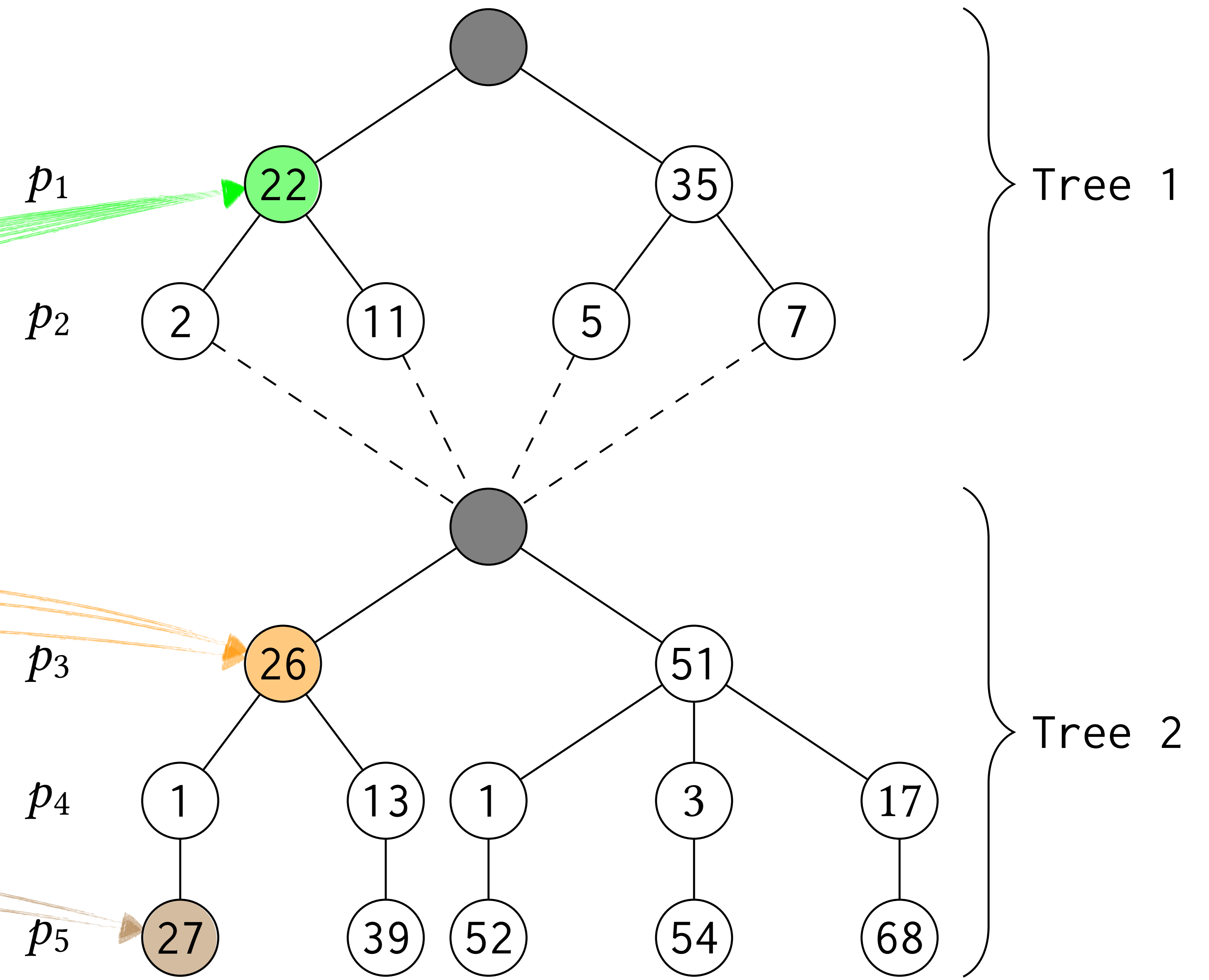
ATF relies on a new *chain-of-trees* search space structure & *parameter constraints*:

$p_1 := (n_1,$	{22, 35},	-)
$p_2 := (n_2,$	{2, 5, 7, 11},	divides(n_1))
$p_3 := (n_3,$	{26, 51},	-)
$p_4 := (n_4,$	{1, 3, 13, 17},	divides(n_3))
$p_5 := (n_5,$	{27, 39, 52, 54, 68},	equals($n_3 + n_4$))

Example Parameters

22	22	22	22	22	22
2	2	2	2	2	11
26	26	51	51	51	26
1	13	1	3	17	1
27	39	52	54	68	27

Traditional Search Space



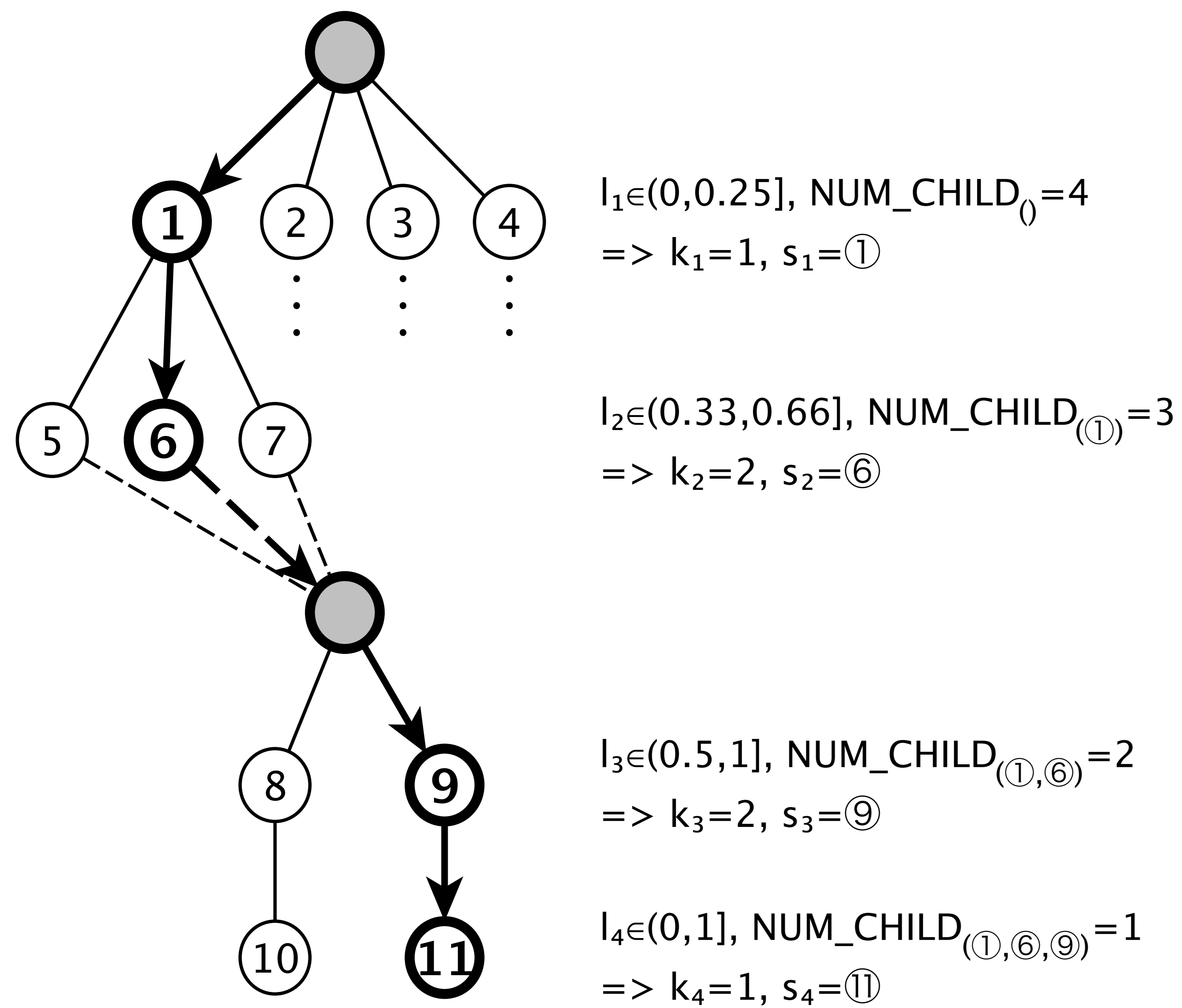
ATF's Search Space

ATF's search space structure avoids memory-intensive redundancies

3. Exploration

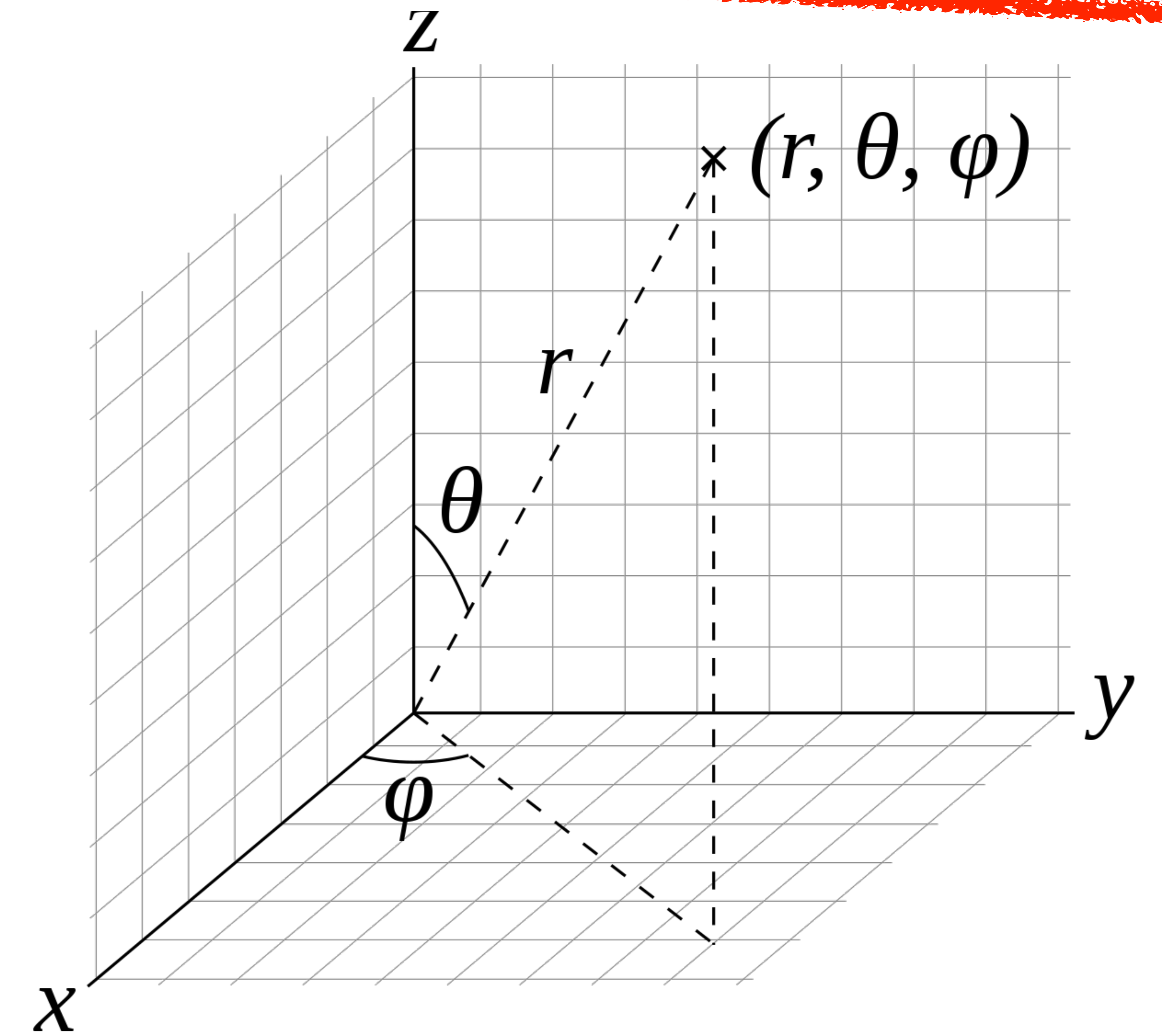
Basic Idea

ATF exploits its new *chain-of-trees* search space structure for a multi-dimensional search:



ATF's Search Space

mapping



Coordinate Space

ATF's search space structure enables reducing the complexity of exploration to exploring a Coordinate Spaces

Experimental Evaluation

Highlights

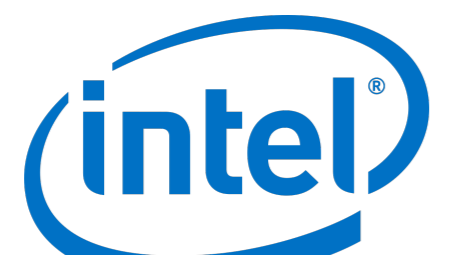
ATF is able to auto-tune important applications for **CPU & GPU** on **real-world data sets** to high performance:


Stencil

ATF is able to auto-tune the CONV implementation in [1] to:

>40x higher performance than CONV+CLTune on CPU

>10⁴x higher performance than CONV+CLTune on GPU

>3x higher performance than Intel MKL-DNN on CPU 

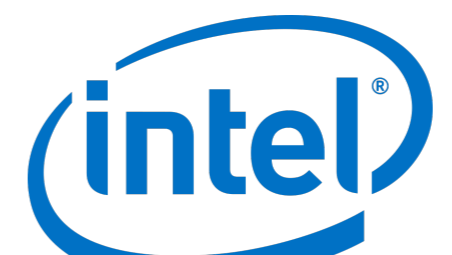
>15x higher performance than NVIDIA cuDNN on GPU 


Linear Algebra

ATF is able to auto-tune the GEMM implementation in [1] to:

>2x higher performance than GEMM+CLTune on CPU

>120x higher performance than GEMM+CLTune on GPU

>2x higher performance than Intel MKL on CPU 

>2x higher performance than NVIDIA cuBLAS on GPU 

Quantum Chemistry

ATF is able to auto-tune the CCSD(T) implementation in [1] to:

>2x higher performance than TensorComprehensions on GPU

CLTune **fails!**
(too high search space generation time)

Data Mining

ATF is able to auto-tune the PRL implementation in [1] to:

>1.66x higher performance than PRL+CLTune on CPU

>1.07x higher performance than PRL+CLTune on GPU

OpenTuner fails for all applications because of a too high amount of invalid configurations within its search space

[1] Rasch, Schulze, Gorlatch. "Generating Portable High-Performance Code via Multi-Dimensional Homomorphisms", PACT'19

ATF — User Interface

ATF's user interface is focus of our previous work [1]:

```
#atf::tp name      NUM_WG_1
      range      interval<int>( 1, N_1 )

#atf::tp name      NUM_WI_1
      range      interval<int>( 1, N_1 )

// ...

#atf::tp name      LM_SIZE_1
      range      interval<int>( 1, N_1 )
      constraint LM_SIZE_1 <= N_1

#atf::tp name      PM_SIZE_1
      range      interval<int>( 1, N_1 )
      constraint PM_SIZE_1 <= LM_SIZE_1

// ...

// OpenCL kernel code
```

**Auto-Tuning in ATF via
easy-to-use Tuning Directives [1]**

(ATF is also available as programming library in C++ [2] / Python (WIP)— for online auto-tuning)

[1] Rasch, Gorlatch. "ATF: A Generic, Directive-Based Auto-Tuning Framework", Concurrency and Computation: Practice and Experience, 2019

[2] Rasch, Haidl, Gorlatch. "ATF: A Generic Auto-Tuning Framework", HPCCC, 2017

Questions?



Ari Rasch
a.rasch@wwu.de