# ATF: A Generic Auto-Tuning Framework

Ari Rasch, Michael Haidl and Sergei Gorlatch

University of Münster, Germany

# Auto-Tuning

**What is Auto-Tuning?**

Approach for automatically optimizing programs:
values of performance-critical (a.k.a. *tuning*) *parameters*, e.g., the
number of threads, are automatically chosen by using a
search technique.

**Why is Auto-Tuning useful?**

- Manually choosing tuning parameter values is hard.

- Optimal values of tuning parameters (usually) differ over target devices
(CPU, GPU, etc.).

# Example:
# SAXPY in OpenCL

- SAXPY is a BLAS routine.

- It takes as its input a scalar a, and two input vectors x and y; it computes:

$$y[i] = a * x[i] + y[i]$$

SAXPY in OpenCL:

- Each thread (a.k.a. work-item) performs computation on chunk of WPT-many elements.

- WPT (Work per Thread) is tuning parameter of the SAXPY kernel.

- Threads are grouped in work-groups.

- Work-group size (a.k.a local size LS) is further tuning parameter of SAXPY kernel.

```
1   __kernel void saxpy( const              int      N,
2                        const              float    a,
3                        const __global float*    x,
4                              __global float*    y
5                      )
6   {
7     for( int w = 0; w < WPT; ++w ) {
8       const int index = w * get_global_size(0)
                          + get_global_id(0);
9
10      y[ index ] += a * x[ index ];
11    }
12  }
```
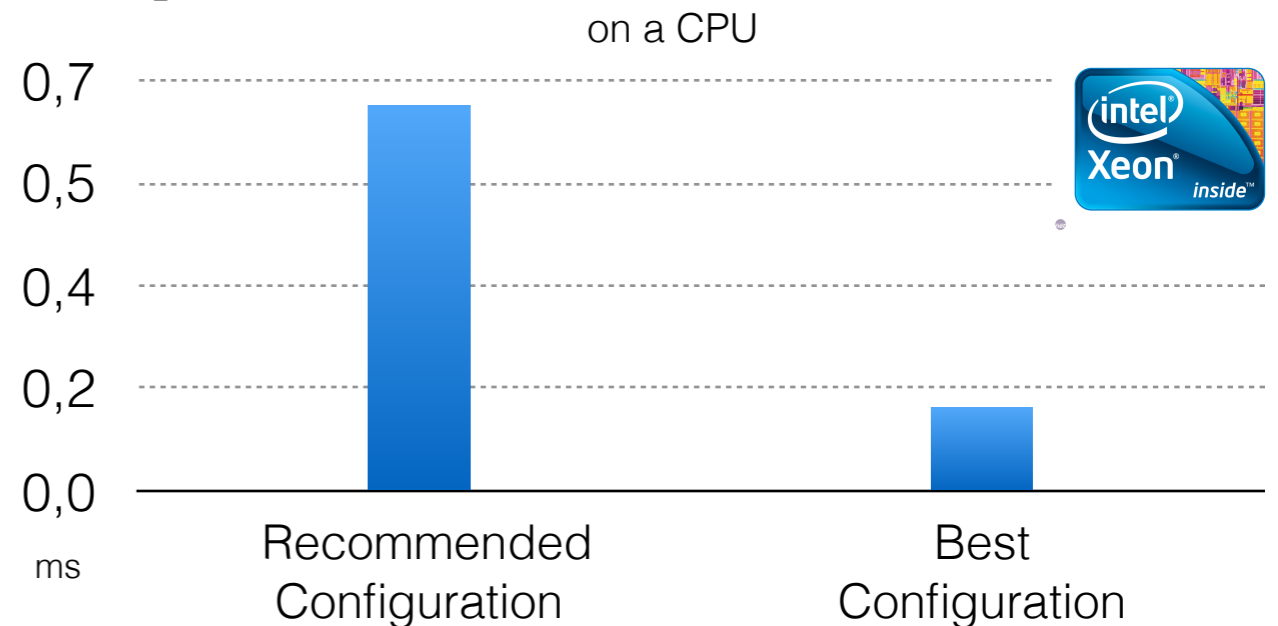
Simplified SAXPY kernel of the auto-tunable CLBlast library
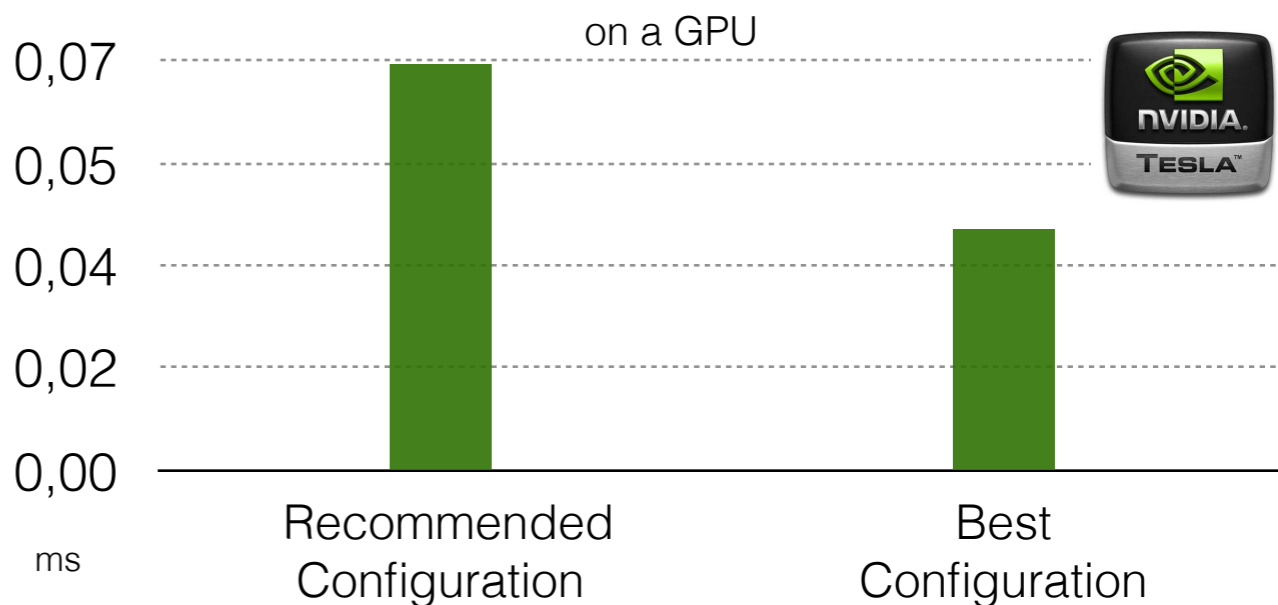
# Example:
# SAXPY in OpenCL

**<u>Manually choosing tuning parameter values is hard:</u>**

- <u>Intel's recommendation for CPU:</u>

  - One work-group per each of CPU's cores;
  - Local size should be 8 (or at least a multiple of 8) enabling SIMD vectorization.

- <u>NVIDIA's recommendation for GPU:</u>

  - As many threads as possible to exploit GPU's massive parallelism;
  - Local size should be 32 (or at least a multiple of 32) to efficiently utilize GPU's *Warps.*

⇒ **Vendor recommendations don't lead to best values!**
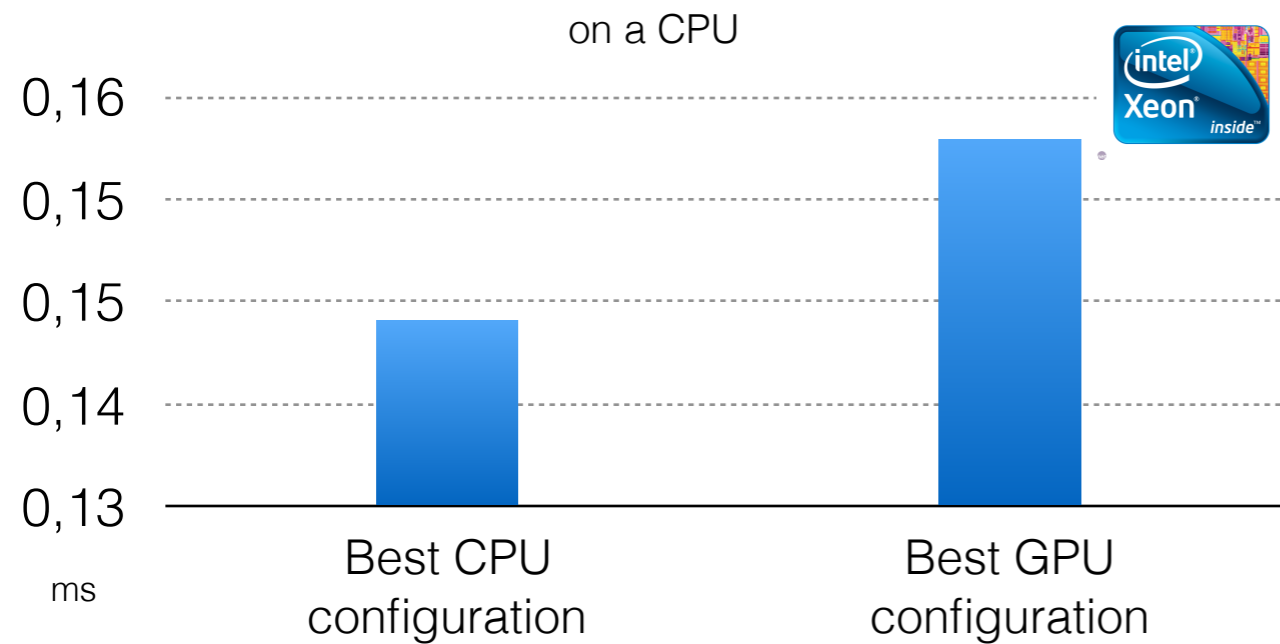
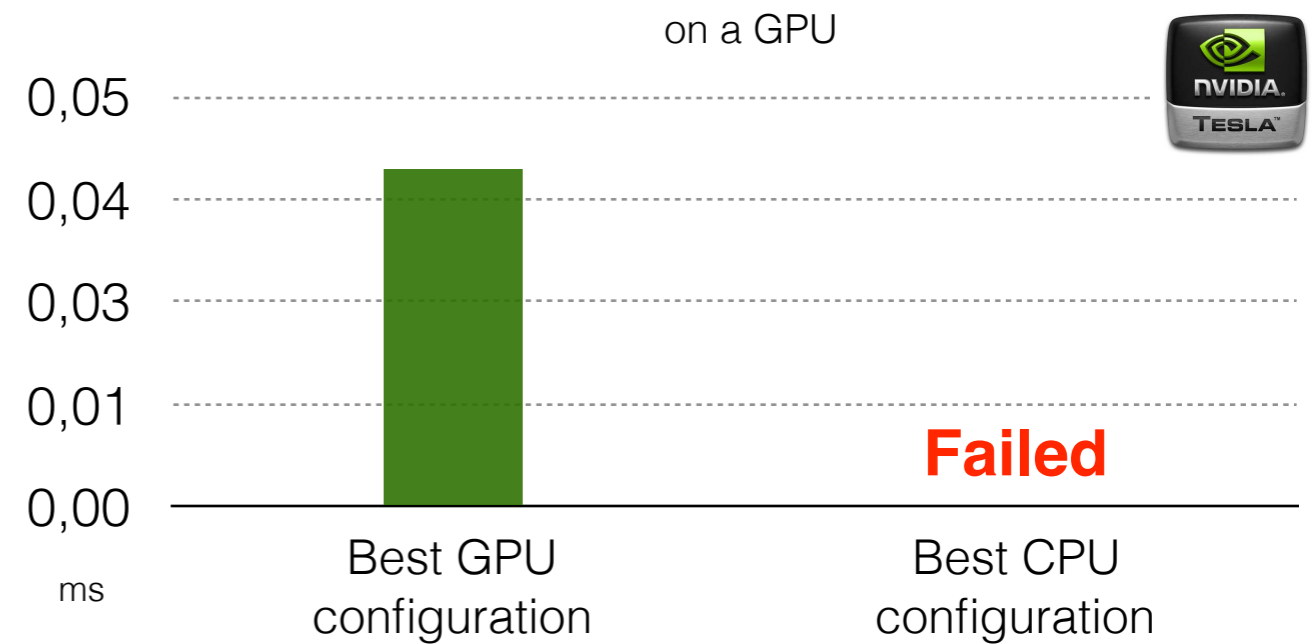on a CPU



| N = 400,000 | Recommended | Best (exhaustive) |
|---|---|---|
| **WPT** | 6250 | 10 |
| **LS** | 8 | 2500 |

on a GPU



| N = 400,000 | Recommended | Best (exhaustive) |
|---|---|---|
| **WPT** | 1 | 4 |
| **LS** | 32 | 160 |

# Example:
# SAXPY in OpenCL

**<u>Optimal values of tuning parameters differ over target devices:</u>**



Configuration with best performance on GPU has not best performance on CPU.

Best configuration for CPU fails on the GPU (unsupported work-group size).

**⇒ Tuning parameter values have to be chosen specifically for each device!**

# Weaknesses of State-of-the-Art Approaches

**Domain-specific approaches:**

- ATLAS (linear algebra), PATUS (stencil), MILEPOST (compiler optimizations), …
  → cannot be used for applications from other domains.

**Generic approaches:**

- OpenTuner: No support for auto-tuning parameters with interdependencies
  (e.g., a parameter has to divide another parameter).

- CLTune: Allows parameters with interdependencies but:

  1. restricted to: 1) only OpenCL, 2) only auto-tuning for runtime performance;

  2. only suitable for auto-tuning small parameter ranges usually not covering
     parameters' entire ranges (search space generation is time intensive).

# The Auto-Tuning Framework (ATF)

ATF combines the following advantages over state-of-the-art auto-tuning approaches:

1. ATF is **generic** regarding:

   - programming language;

   - application domain;

   - tuning objective;

   - search technique.

2. ATF allows **dependencies between tuning parameters**.

3. ATF allows **significantly larger tuning parameter ranges**.

4. ATF is **arguably simpler** to use.

# Illustration of ATF

- Illustration of ATF by a simple example: auto-tuning the OpenCL SAXPY kernel.

- The ATF user has to implement a C++ program using ATF's C++ API and perform the following three steps:

## 1. Step: Define the search space

- ATF search spaces are defined using tuning parameters, here:

  - WPT: a size_t parameter in [1,N] that divides N

  - LS: a size_t parameter in [1,N] that divides N/WPT

```cpp
1   int main()
2   {
3     std::string saxpy_kernel = /* SAXPY kernel's code as string */;
4     int          N           = /* fixed user-defined input size */;
5
6     auto WPT = atf::tp( "WPT",
7                         atf::interval<size_t>(1,N),
8                         atf::divides( N )
9                       );
10    auto LS  = atf::tp( "LS",
11                        atf::interval<size_t>(1,N),
12                        atf::divides( N/WPT )
13                      );
14
15    auto cf_saxpy = atf::cf::ocl(
16                      { "NVIDIA", "Tesla K20c" },
17                      saxpy_kernel,
18                      inputs( atf::scalar<int>(N),   // N
19                              atf::scalar<float>(),  // a
20                              atf::buffer<float>(N), // x
21                              atf::buffer<float>(N), // y
22                            )
23                      atf::glb_size( N/WPT ), atf::lcl_size( LS )
24                    );
25
26    auto best_config = atf::annealing( atf::duration<minutes>(10) )
27                         ( WPT, LS )
28                         ( cf_saxpy );
29  }
```

ATF program for auto-tuning SAXPY

# Illustration of ATF

## 2. Step: Implement a cost function

- Cost function: Configuration → Cost (e.g., runtime)

- Here, we use ATF's pre-implemented cost function for auto-tuning OpenCL in terms of runtime:

  - it is initialized with: i) target device, ii) kernel's code, iii) kernels' input arguments, iv) global/ local size;

  - tuning parameters are replaced by values according to the input configuration;

  - it returns kernel's runtime as cost.

## 3. Step: Explore the search space

- The search space is explored by choosing a search technique and passing to it:

  1) an abort condition

  2) the tuning parameters

  3) the cost function

- The result is the best found configuration.

```
1   int main()
2   {
3     std::string saxpy_kernel = /* SAXPY kernel's code as string */;
4     int         N            = /* fixed user-defined input size */;
5
6     auto WPT = atf::tp( "WPT",
7                         atf::interval<size_t>(1,N),
8                         atf::divides( N )
9                       );
10    auto LS  = atf::tp( "LS",
11                        atf::interval<size_t>(1,N),
12                        atf::divides( N/WPT )
13                      );
14
15    auto cf_saxpy = atf::cf::ocl(
16                      { "NVIDIA", "Tesla K20c" },
17                      saxpy_kernel,
18                      inputs( atf::scalar<int>(N),    // N
19                              atf::scalar<float>(),   // a
20                              atf::buffer<float>(N),  // x
21                              atf::buffer<float>(N),  // y
22                            )
23                      atf::glb_size( N/WPT ), atf::lcl_size( LS )
24                    );
25
26    auto best_config = atf::annealing( atf::duration<minutes>(10) )
27                       ( WPT, LS )
28                       ( cf_saxpy );
29  }
```

ATF program for auto-tuning SAXPY

# Detailed Discussion of ATF

## 1. Step: Define the search space

General form of a tuning parameter:

```
atf::tp( /* name        */,
         /* range       */,
         /* constraints */
);
```

Used to access parameter's value in a configuration , e.g., best_config["LS"].

Boolean functions to filter the parameter's range; may contain tuning parameters to express interdependencies, e.g.: atf::divides( N/WPT ) (→ [&](auto LS){ N/WPT % LS == 0} )

Can be either an:

1) atf::interval<T>(begin, end, step_size=1, generator=id), where generator: T → U

2) atf::set( val_1, … , val_n)

# Detailed Discussion of ATF

## 2. Step: Implement a cost function

General form of a cost function:

```
T cost_function( atf::configuration config )
{
  // ...
}
```

- Input: a parameter configuration

- Output: Element of type T (e.g., size_t) for which < is defined

- ATF interprets output as cost to minimize (e.g., runtime).

- Multi-Objective Tuning: e.g., runtime & energy → T=std::array<size_t,2> with < as lexicographical order

- ATF provides three pre-implemented cost functions, for:

  1. OpenCL,

  2. CUDA,

  3. arbitrary programming languages that are not OpenCL or CUDA.

```
atf::cf::ocl(
  {"NVIDIA", "Tesla K20c"},
  saxpy,
  inputs( atf::scalar<int>(N),    // N
          atf::scalar<float>(),   // a
          atf::buffer<float>(N),  // x
          atf::buffer<float>(N),  // y
        )
  atf::glb_size( N/WPT ),
  atf::lcl_size( LS )
);
```

OpenCL

```
atf::cf::cuda(
  {"Tesla K20c"},
  saxpy_cuda,
  inputs( atf::scalar<int>(N),    // N
          atf::scalar<float>(),   // a
          atf::buffer<float>(N),  // x
          atf::buffer<float>(N),  // y
        )
  atf::grd_size( (N/WPT)/LS ),
  atf::blk_size( LS )
);
```

CUDA

```
atf::cf::gcf(
  /* path to source file   */,
  /* path to compile script */,
  /* path to run script    */,
  /* path to cost file     */
);
```

Generic Cost Function

# Detailed Discussion of ATF

## 3. Step: Explore the search space

General schema for exploring the search space:

```
atf::/∗ search technique ∗/( /∗ abort condition   ∗/ )
                           ( /∗ tuning parameters ∗/ )
                           ( /∗ cost function     ∗/ );
```

**ATF provides three pre-implemented search techniques:**

1. Exhaustive: finds provably best configuration, but probably at the cost of a long search time;

2. Simulated Annealing: effective for auto-tuning OpenCL/CUDA for large search spaces;

3. OpenTuner: combines various techniques to yield a good tuning result on average for arbitrary applications.

**ATF provides various abort conditions, e.g.:**

- duration<D>(t): stops tuning after time interval t; D is an std::chrono::duration (sec, min, etc.)

- cost(c): stops tuning when a configuration with a cost lower or equal than c has been found;

- speedup<D>(s,t): stops tuning when within last time interval t cost could not be lowered by a factor >=s;

- …

# Comparison: ATF vs. CLTune

- <u>We show:</u> even though ATF is a generic approach, it works better for OpenCL than CLTune which is specifically designed for OpenCL.

- We use the example of auto-tuning SAXPY.

```cpp
int main()
{
  const std::string saxpy = /* path to kernel of Listing 1   */;
  const size_t      N      = /* fixed user-defined input size */;

  cltune::Tuner tuner(1,0);
  auto id = tuner.AddKernel(saxpy, "saxpy", {N}, {1} );

  float a;
  auto  x = std::vector<float>(N);
  auto  y = std::vector<float>(N);

  const auto random_seed =
    std::chrono::system_clock::now().time_since_epoch().count();
  std::default_random_engine
    generator( static_cast<unsigned int>(random_seed) );
  std::uniform_real_distribution<float> distribution(-2.0f,2.0f);

  a = distribution(generator);
  for (auto &item: x) { item = distribution(generator); }
  for (auto &item: y) { item = distribution(generator); }

  tuner.AddArgumentScalar( N );
  tuner.AddArgumentScalar( a );
  tuner.AddArgumentInput( x );
  tuner.AddArgumentOutput( y );

  auto range = std::vector<size_t>( N );
  for( size_t i = 0; i < N ; ++i )
    range[ i ] = i;
  tuner.AddParameter( id, "LS" , range );
  tuner.AddParameter( id, "WPT", range );

  auto DividesN = []( std::vector<size_t> v )
                  {
                    return  N % v[0] == 0;
                  };
  auto DividesNDivWPT = []( std::vector<size_t> v )
                  {
                    return ( N / v[0] ) % v[1] == 0;
                  };

  tuner.AddConstraint( id, DividesN       , {"WPT"}       );
  tuner.AddConstraint( id, DividesNDivWPT, {"WPT", "LS"} );

  tuner.DivGlobalSize(id, {"WPT" } );
  tuner.MulLocalSize(id, {"LS"} );

  tuner.UseAnnealing( 1.0f/2048.0f , 4.0 );
  tuner.Tune();
  const auto parameters = tuner.GetBestResult();
}
```

# Comparison: ATF vs. CLTune

```
1   int main()
2   {
3     std::string saxpy_kernel = /* path to kernel of Listing 1   */;
4     int         N            = /* fixed user-defined input size */;
5
6     auto WPT = atf::tp( "WPT",
7                         atf::interval<size_t>(1,N),
8                         atf::divides( N )
9                       );
10    auto LS  = atf::tp( "LS",
11                        atf::interval<size_t>(1,N),
12                        atf::divides( N/WPT )
13                      );
14
15    auto cf_saxpy = atf::cf::ocl(
16                      { "NVIDIA", "Tesla K20c" },
17                      saxpy_kernel,
18                      inputs( atf::scalar<int>(N),    // N
19                              atf::scalar<float>(),   // a
20                              atf::buffer<float>(N), // x
21                              atf::buffer<float>(N), // y
22                            )
23                      atf::glb_size( N/WPT ), atf::lcl_size( LS )
24                    );
25
26    auto best_config = atf::annealing( atf::duration<minutes>(10) )
27                                      ( WPT, LS )
28                                      ( cf_saxpy );
29  }
```

```
1   int main()
2   {
3     const std::string saxpy = /* path to kernel of Listing 1   */;
4     const size_t      N     = /* fixed user-defined input size */;
5
6     cltune::Tuner tuner(1,0);
7     auto id = tuner.AddKernel(saxpy, "saxpy", {N}, {1} );
8
9     float a;
10    auto  x = std::vector<float>(N);
11    auto  y = std::vector<float>(N);
12
13    const auto random_seed =
        std::chrono::system_clock::now().time_since_epoch().count();
14    std::default_random_engine
        generator( static_cast<unsigned int>(random_seed) );
15    std::uniform_real_distribution<float> distribution(-2.0f,2.0f);
16
17    a = distribution(generator);
18    for (auto &item: x) { item = distribution(generator); }
19    for (auto &item: y) { item = distribution(generator); }
20
21    tuner.AddArgumentScalar( N );
22    tuner.AddArgumentScalar( a );
23    tuner.AddArgumentInput( x );
24    tuner.AddArgumentOutput( y );
25
26    auto range = std::vector<size_t>( N );
27    for( size_t i = 0; i < N ; ++i )
28      range[ i ] = i;
29    tuner.AddParameter( id, "LS" , range );
30    tuner.AddParameter( id, "WPT", range );
31
32    auto DividesN = []( std::vector<size_t> v )
                    {
                      return  N % v[0] == 0;
                    };
33    auto DividesNDivWPT = []( std::vector<size_t> v )
                    {
                      return ( N / v[0] ) % v[1] == 0;
                    };
34
35    tuner.AddConstraint( id, DividesN       , {"WPT"}        );
36    tuner.AddConstraint( id, DividesNDivWPT, {"WPT", "LS"} );
37
38    tuner.DivGlobalSize(id, {"WPT" } );
39    tuner.MulLocalSize(id, {"LS"} );
40
41    tuner.UseAnnealing( 1.0f/2048.0f , 4.0 );
42    tuner.Tune();
43    const auto parameters = tuner.GetBestResult();
44  }
```

# Comparison: ATF vs. CLTune

```
1   int main()
2   {
3     std::string saxpy_kernel = /* path to kernel of Listing 1   */;
4     int         N            = /* fixed user-defined input size */;
5
6     auto WPT = atf::tp( "WPT",
7                         atf::interval<size_t>(1,N),
8                         atf::divides( N )
9                       );
10    auto LS  = atf::tp( "LS",
11                        atf::interval<size_t>(1,N),
12                        atf::divides( N/WPT )
13                      );
14
15    auto cf_saxpy = atf::cf::ocl(
16                       { "NVIDIA", "Tesla K20c" },
17                       saxpy_kernel,
18                       inputs( atf::scalar<int>(N),   // N
19                               atf::scalar<float>(),  // a
20                               atf::buffer<float>(N), // x
21                               atf::buffer<float>(N), // y
22                             )
23                       atf::glb_size( N/WPT ), atf::lcl_size( LS )
24                     );
25
26    auto best_config = atf::annealing( atf::duration<minutes>(10) )
27                                      ( WPT, LS )
28                                      ( cf_saxpy );
29  }
```

## ATF allows easier expressing parameter dependencies

```
1   int main()
2   {
3     const std::string saxpy = /* path to kernel of Listing 1   */;
4     const size_t       N     = /* fixed user-defined input size */;
5
6     cltune::Tuner tuner(1,0);
7     auto id = tuner.AddKernel(saxpy, "saxpy", {N}, {1} );
8
9     float a;
10    auto  x = std::vector<float>(N);
11    auto  y = std::vector<float>(N);
12
13    const auto random_seed =
        std::chrono::system_clock::now().time_since_epoch().count();
14    std::default_random_engine
        generator( static_cast<unsigned int>(random_seed) );
15    std::uniform_real_distribution<float> distribution(-2.0f,2.0f);
16
17    a = distribution(generator);
18    for (auto &item: x) { item = distribution(generator); }
19    for (auto &item: y) { item = distribution(generator); }
20
21    tuner.AddArgumentScalar( N );
22    tuner.AddArgumentScalar( a );
23    tuner.AddArgumentInput( x );
24    tuner.AddArgumentOutput( y );
25
26    auto range = std::vector<size_t>( N );
27    for( size_t i = 0; i < N ; ++i )
28      range[ i ] = i;
29    tuner.AddParameter( id, "LS" , range );
30    tuner.AddParameter( id, "WPT", range );
31
32    auto DividesN = []( std::vector<size_t> v )
                      {
                        return  N % v[0] == 0;
                      };
33    auto DividesNDivWPT = []( std::vector<size_t> v )
                           {
                             return ( N / v[0] ) % v[1] == 0;
                           };
34
35    tuner.AddConstraint( id, DividesN       , {"WPT"}       );
36    tuner.AddConstraint( id, DividesNDivWPT, {"WPT", "LS"} );
37
38    tuner.DivGlobalSize(id, {"WPT" } );
39    tuner.MulLocalSize(id, {"LS"} );
40
41    tuner.UseAnnealing( 1.0f/2048.0f , 4.0 );
42    tuner.Tune();
43    const auto parameters = tuner.GetBestResult();
44  }
```

# Comparison: ATF vs. CLTune

```
1   int main()
2   {
3     std::string saxpy_kernel = /* path to kernel of Listing 1   */;
4     int          N           = /* fixed user-defined input size */;
5
6     auto WPT = atf::tp( "WPT",
7                         atf::interval<size_t>(1,N),
8                         atf::divides( N )
9                       );
10    auto LS  = atf::tp( "LS",
11                        atf::interval<size_t>(1,N),
12                        atf::divides( N/WPT )
13                      );
14
15    auto cf_saxpy = atf::cf::ocl(
16                        { "NVIDIA", "Tesla K20c" },
17                        saxpy_kernel,
18                        inputs( atf::scalar<int>(N),    // N
19                                atf::scalar<float>(),   // a
20                                atf::buffer<float>(N),  // x
21                                atf::buffer<float>(N),  // y
22                              )
23                        atf::glb_size( N/WPT ), atf::lcl_size( LS )
24                    );
25
26    auto best_config = atf::annealing( atf::duration<minutes>(10) )
27                                     ( WPT, LS )
28                                     ( cf_saxpy );
29  }
```

```
1   int main()
2   {
3     const std::string saxpy = /* path to kernel of Listing 1   */;
4     const size_t       N     = /* fixed user-defined input size */;
5
6     cltune::Tuner tuner(1,0);
7     auto id = tuner.AddKernel(saxpy, "saxpy", {N}, {1} );
8
9     float a;
10    auto  x = std::vector<float>(N);
11    auto  y = std::vector<float>(N);
12
13    const auto random_seed =
        std::chrono::system_clock::now().time_since_epoch().count();
14    std::default_random_engine
        generator( static_cast<unsigned int>(random_seed) );
15    std::uniform_real_distribution<float> distribution(-2.0f,2.0f);
16
17    a = distribution(generator);
18    for (auto &item: x) { item = distribution(generator); }
19    for (auto &item: y) { item = distribution(generator); }
20
21    tuner.AddArgumentScalar( N );
22    tuner.AddArgumentScalar( a );
23    tuner.AddArgumentInput( x );
24    tuner.AddArgumentOutput( y );
25
26    auto range = std::vector<size_t>( N );
27    for( size_t i = 0; i < N ; ++i )
28      range[ i ] = i;
29    tuner.AddParameter( id, "LS" , range );
30    tuner.AddParameter( id, "WPT", range );
31
32    auto DividesN = []( std::vector<size_t> v )
                     {
                         return  N % v[0] == 0;
                     };
33    auto DividesNDivWPT = []( std::vector<size_t> v )
                           {
                               return ( N / v[0] ) % v[1] == 0;
                           };
34
35    tuner.AddConstraint( id, DividesN       , {"WPT"}      );
36    tuner.AddConstraint( id, DividesNDivWPT, {"WPT", "LS"} );
37
38    tuner.DivGlobalSize(id, {"WPT" } );
39    tuner.MulLocalSize(id, {"LS"} );
40
41    tuner.UseAnnealing( 1.0f/2048.0f , 4.0 );
42    tuner.Tune();
43    const auto parameters = tuner.GetBestResult();
44  }
```

## ATF allows easier setting the global/local size

16 / 24

# Comparison: ATF vs. CLTune

```cpp
1   int main()
2   {
3     std::string saxpy_kernel = /* path to kernel of Listing 1   */;
4     int          N           = /* fixed user-defined input size */;
5
6     auto WPT = atf::tp( "WPT",
7                         atf::interval<size_t>(1,N),
8                         atf::divides( N )
9                       );
10    auto LS  = atf::tp( "LS",
11                        atf::interval<size_t>(1,N),
12                        atf::divides( N/WPT )
13                      );
14
15    auto cf_saxpy = atf::cf::ocl(
16                    { "NVIDIA", "Tesla K20c" },
17                    saxpy_kernel,
18                    inputs( atf::scalar<int>(N),    // N
19                            atf::scalar<float>(),   // a
20                            atf::buffer<float>(N),  // x
21                            atf::buffer<float>(N),  // y
22                          )
23                    atf::glb_size( N/WPT ), atf::lcl_size( LS )
24                  );
25
26    auto best_config = atf::annealing( atf::duration<minutes>(10) )
27                                     ( WPT, LS )
28                                     ( cf_saxpy );
29  }
```

## ATF allows a broader range of global/local sizes

```cpp
1   int main()
2   {
3     const std::string saxpy = /* path to kernel of Listing 1   */;
4     const size_t       N     = /* fixed user-defined input size */;
5
6     cltune::Tuner tuner(1,0);
7     auto id = tuner.AddKernel(saxpy, "saxpy", {N}, {1} );
8
9     float a;
10    auto  x = std::vector<float>(N);
11    auto  y = std::vector<float>(N);
12
13    const auto random_seed =
          std::chrono::system_clock::now().time_since_epoch().count();
14    std::default_random_engine
          generator( static_cast<unsigned int>(random_seed) );
15    std::uniform_real_distribution<float> distribution(-2.0f,2.0f);
16
17    a = distribution(generator);
18    for (auto &item: x) { item = distribution(generator); }
19    for (auto &item: y) { item = distribution(generator); }
20
21    tuner.AddArgumentScalar( N );
22    tuner.AddArgumentScalar( a );
23    tuner.AddArgumentInput( x );
24    tuner.AddArgumentOutput( y );
25
26    auto range = std::vector<size_t>( N );
27    for( size_t i = 0; i < N ; ++i )
28      range[ i ] = i;
29    tuner.AddParameter( id, "LS" , range );
30    tuner.AddParameter( id, "WPT", range );
31
32    auto DividesN = []( std::vector<size_t> v )
                      {
                        return  N % v[0] == 0;
                      };
33    auto DividesNDivWPT = []( std::vector<size_t> v )
                           {
                             return ( N / v[0] ) % v[1] == 0;
                           };
34
35    tuner.AddConstraint( id, DividesN      , {"WPT"}     );
36    tuner.AddConstraint( id, DividesNDivWPT, {"WPT", "LS"} );
37
38    tuner.DivGlobalSize(id, {"WPT" } );
39    tuner.MulLocalSize(id, {"LS"} );
40
41    tuner.UseAnnealing( 1.0f/2048.0f , 4.0 );
42    tuner.Tune();
43    const auto parameters = tuner.GetBestResult();
44  }
```

# Comparison: ATF vs. CLTune

```cpp
1  int main()
2  {
3    std::string saxpy_kernel = /* path to kernel of Listing 1   */;
4    int         N            = /* fixed user-defined input size */;
5
6    auto WPT = atf::tp( "WPT",
7                        atf::interval<size_t>(1,N),
8                        atf::divides( N )
9                      );
10   auto LS  = atf::tp( "LS",
11                       atf::interval<size_t>(1,N),
12                       atf::divides( N/WPT )
13                     );
14
15   auto cf_saxpy = atf::cf::ocl(
16                     { "NVIDIA", "Tesla K20c" },
17                     saxpy_kernel,
18                     inputs( atf::scalar<int>(N),    // N
19                             atf::scalar<float>(),   // a
20                             atf::buffer<float>(N),  // x
21                             atf::buffer<float>(N),  // y
22                           )
23                     atf::glb_size( N/WPT ), atf::lcl_size( LS )
24                  );
25
26   auto best_config = atf::annealing( atf::duration<minutes>(10) )
27                                    ( WPT, LS )
28                                    ( cf_saxpy );
29 }
```

```cpp
1  int main()
2  {
3    const std::string saxpy = /* path to kernel of Listing 1   */;
4    const size_t      N     = /* fixed user-defined input size */;
5
6    cltune::Tuner tuner(1,0);
7    auto id = tuner.AddKernel(saxpy, "saxpy", {N}, {1} );
8
9    float a;
10   auto  x = std::vector<float>(N);
11   auto  y = std::vector<float>(N);
12
13   const auto random_seed =
        std::chrono::system_clock::now().time_since_epoch().count();
14   std::default_random_engine
        generator( static_cast<unsigned int>(random_seed) );
15   std::uniform_real_distribution<float> distribution(-2.0f,2.0f);
16
17   a = distribution(generator);
18   for (auto &item: x) { item = distribution(generator); }
19   for (auto &item: y) { item = distribution(generator); }
20
21   tuner.AddArgumentScalar( N );
22   tuner.AddArgumentScalar( a );
23   tuner.AddArgumentInput( x );
24   tuner.AddArgumentOutput( y );
25
26   auto range = std::vector<size_t>( N );
27   for( size_t i = 0; i < N ; ++i )
28     range[ i ] = i;
29   tuner.AddParameter( id, "LS" , range );
30   tuner.AddParameter( id, "WPT", range );
31
32   auto DividesN = []( std::vector<size_t> v )
                    {
                      return  N % v[0] == 0;
                    };
33   auto DividesNDivWPT = []( std::vector<size_t> v )
                          {
                            return ( N / v[0] ) % v[1] == 0;
                          };
34
35   tuner.AddConstraint( id, DividesN       , {"WPT"}        );
36   tuner.AddConstraint( id, DividesNDivWPT, {"WPT", "LS"} );
37
38   tuner.DivGlobalSize(id, {"WPT" } );
39   tuner.MulLocalSize(id, {"LS"} );
40
41   tuner.UseAnnealing( 1.0f/2048.0f , 4.0 );
42   tuner.Tune();
43   const auto parameters = tuner.GetBestResult();
44 }
```

## ATF allows easier generating random input data

# Comparison: ATF vs. CLTune

```
1   int main()
2   {
3     std::string saxpy_kernel = /* path to kernel of Listing 1   */;
4     int         N            = /* fixed user-defined input size */;
5
6     auto WPT = atf::tp( "WPT",
7                         atf::interval<size_t>(1,N),
8                         atf::divides( N )
9                       );
10    auto LS  = atf::tp( "LS",
11                        atf::interval<size_t>(1,N),
12                        atf::divides( N/WPT )
13                      );
14
15    auto cf_saxpy = atf::cf::ocl(
16                      { "NVIDIA", "Tesla K20c" },
17                      saxpy_kernel,
18                      inputs( atf::scalar<int>(N),   // N
19                              atf::scalar<float>(),  // a
20                              atf::buffer<float>(N), // x
21                              atf::buffer<float>(N), // y
22                            )
23                      atf::glb_size( N/WPT ), atf::lcl_size( LS )
24                    );
25
26    auto best_config = atf::annealing( atf::duration<minutes>(10) )
27                                     ( WPT, LS )
28                                     ( cf_saxpy );
29  }
```
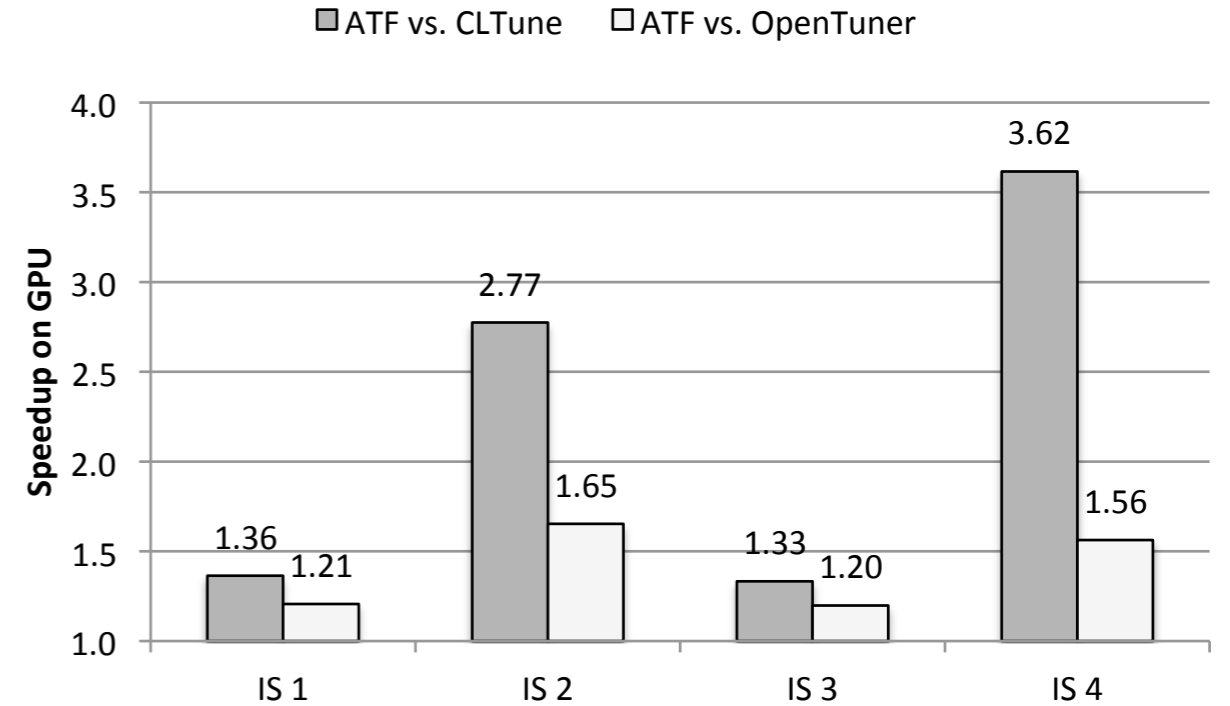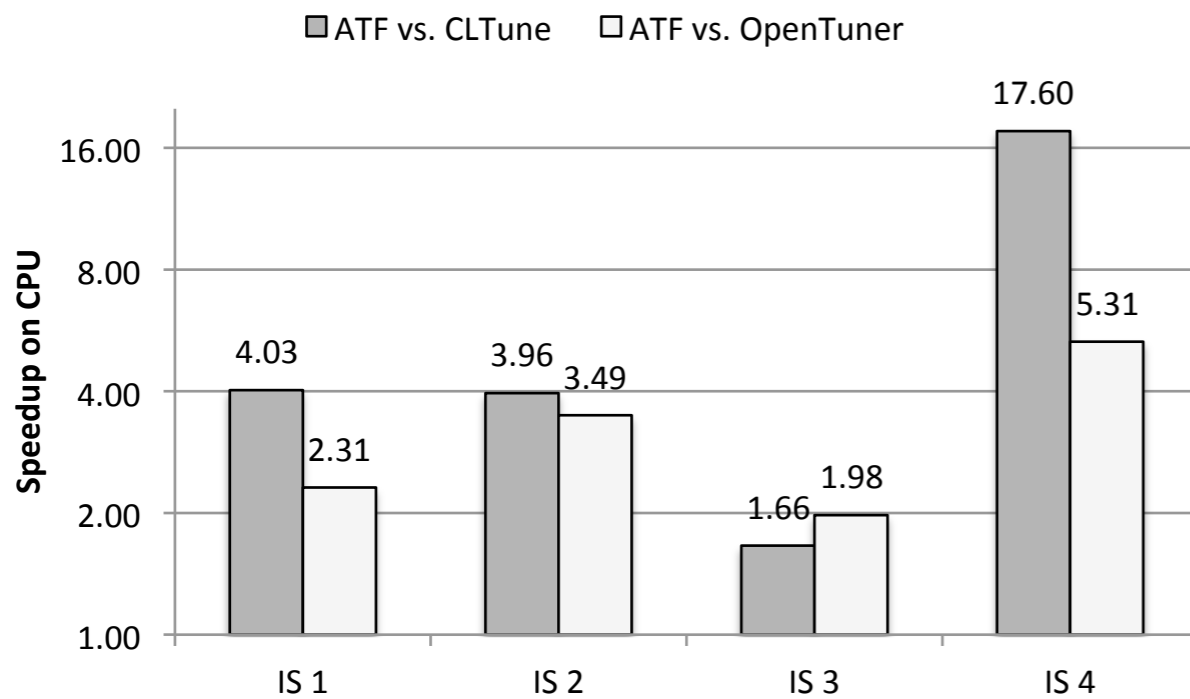
```
1   int main()
2   {
3     const std::string saxpy = /* path to kernel of Listing 1   */;
4     const size_t      N     = /* fixed user-defined input size */;
5
6     cltune::Tuner tuner(1,0);
7     auto id = tuner.AddKernel(saxpy, "saxpy", {N}, {1} );
8
9     float a;
10    auto  x = std::vector<float>(N);
11    auto  y = std::vector<float>(N);
12
13    const auto random_seed =
14      std::chrono::system_clock::now().time_since_epoch().count();
      std::default_random_engine
        generator( static_cast<unsigned int>(random_seed) );
15    std::uniform_real_distribution<float> distribution(-2.0f,2.0f);
16
17    a = distribution(generator);
18    for (auto &item: x) { item = distribution(generator); }
19    for (auto &item: y) { item = distribution(generator); }
20
21    tuner.AddArgumentScalar( N );
22    tuner.AddArgumentScalar( a );
23    tuner.AddArgumentInput( x );
24    tuner.AddArgumentOutput( y );
25
26    auto range = std::vector<size_t>( N );
27    for( size_t i = 0; i < N ; ++i )
28      range[ i ] = i;
29    tuner.AddParameter( id, "LS" , range );
30    tuner.AddParameter( id, "WPT", range );
31
32    auto DividesN = []( std::vector<size_t> v )
                      {
                        return  N % v[0] == 0;
                      };
33    auto DividesNDivWPT = []( std::vector<size_t> v )
                            {
                              return ( N / v[0] ) % v[1] == 0;
                            };
34
35    tuner.AddConstraint( id, DividesN      , {"WPT"}      );
36    tuner.AddConstraint( id, DividesNDivWPT, {"WPT", "LS"} );
37
38    tuner.DivGlobalSize(id, {"WPT" } );
39    tuner.MulLocalSize(id, {"LS"} );
40
41    tuner.UseAnnealing( 1.0f/2048.0f , 4.0 );
42    tuner.Tune();
43    const auto parameters = tuner.GetBestResult();
44  }
```

## ATF supports more abort conditions

# Experimental Results

- <u>We demonstrate:</u> ATF provides better tuning results for GEMM (GEneral Matrix Multiplication) than OpenTuner and CLTune.

- As concrete GEMM implementation, we use the OpenCL kernel XgemmDirect which is part of the CLBlast library that uses CLTune for auto-tuning.

- XgemmDirect is used for accelerating important applications, e.g., the state-of-the-art deep learning framework Caffe [Jia et al, 2014].

- XgemmDirect has 10 tuning parameter, e.g., tile size WGD and loop unrolling factor KWID.

- Tuning parameters have various interdependencies (16 in total), e.g, KWID has to divide WGD.

- We study four pairs of matrix input sizes (IS) with significance in deep learning:

  - <u>IS 1:</u> 20×1     and  1×576

  - <u>IS 2:</u> 20×25   and 25×576

  - <u>IS 3:</u> 50×1     and   1×64

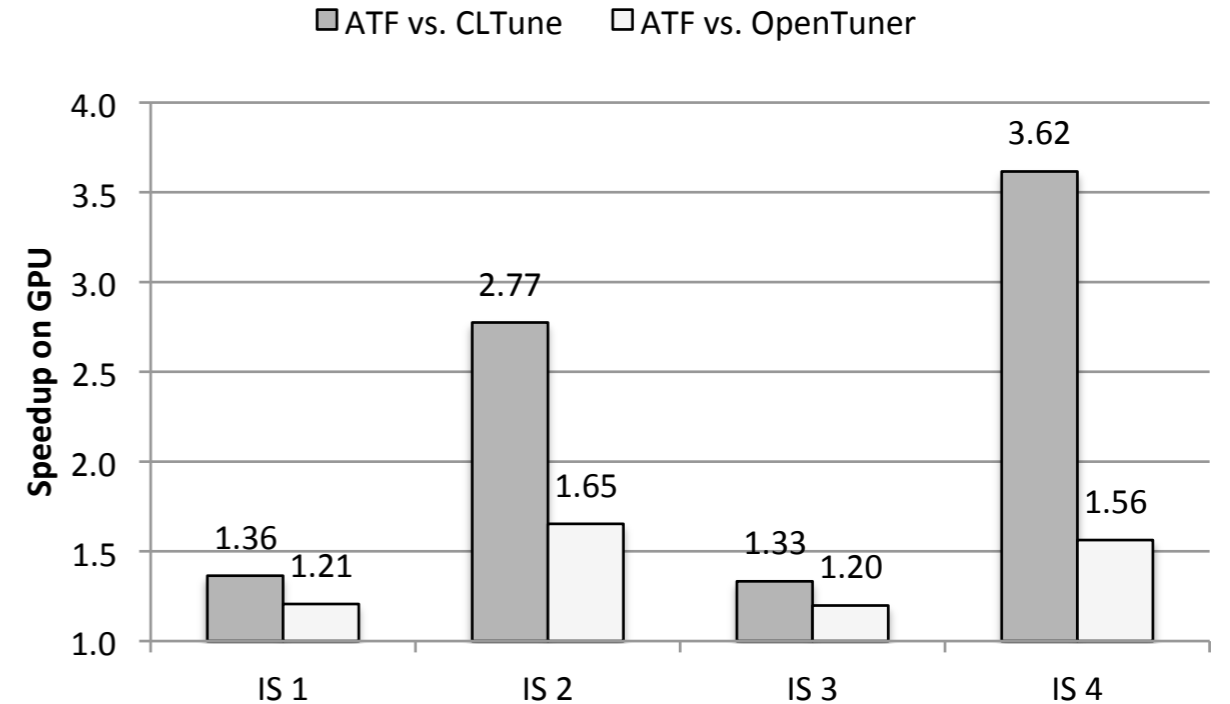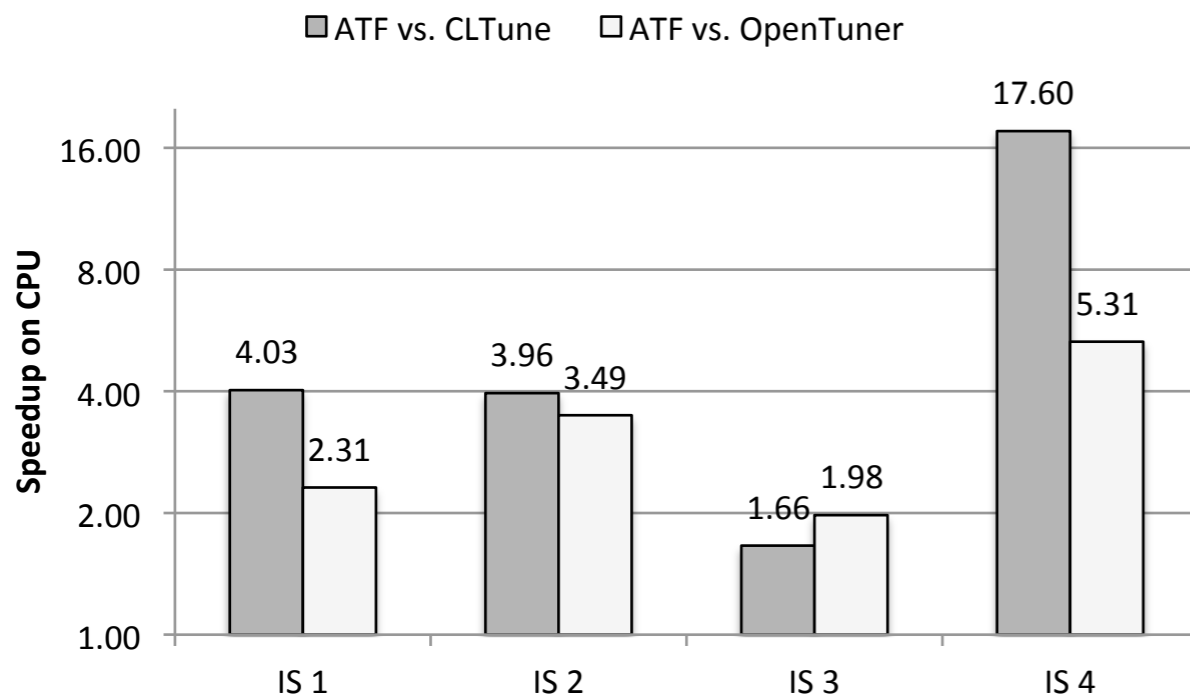  - <u>IS 4:</u> 10×64   and 64x500

# Experimental Results



Speedup of XgemmDirect auto-tuned by ATF compared to CLTune and OpenTuner.

ATF vs CLTune:

- Speedups of up to **17.60x** on CPU and **3.62x** on GPU:

  - CLTune uses artificially limited tuning parameter ranges to shorten its time-intensive search space generation.

  - Even though, the limitations are performed by an expert optimal solutions are missed.

  - ATF uses unlimited ranges, comprising good values also for the special matrix sizes used in deep learning.

- Removing CLTune's artificial limitations causes high search space time: aborted CLTune after 3 hours; ATF requires less than 1 second → ATF filters parameter ranges while CLTune filters the (large) search space.

# Experimental Results



Speedup of XgemmDirect auto-tuned by ATF compared to XgemmDirect auto-tuned by CLTune and OpenTuner.
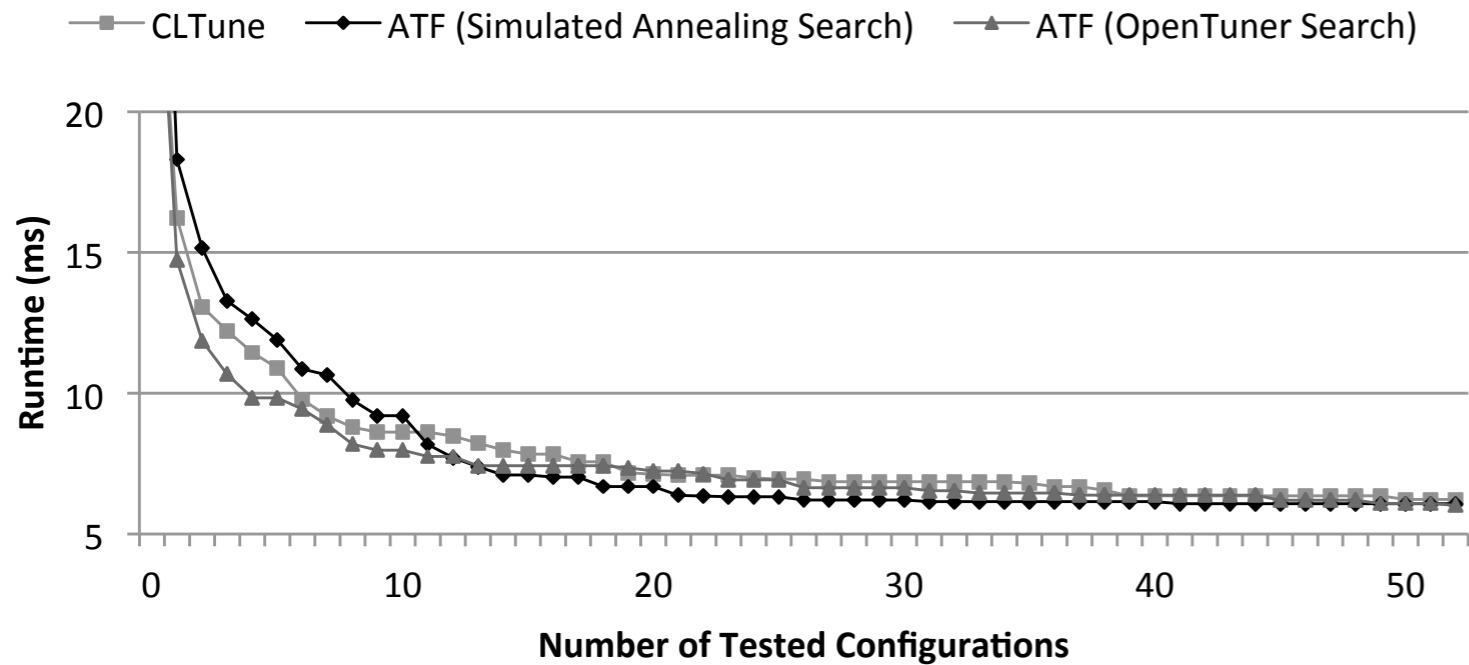
ATF vs OpenTuner:

- Speedups of up to **5.31x** on CPU and **1.65x** on GPU:

  - OpenTuner uses unconstrained search space and can't find valid configurations.

  - Search space size for IS 4: $10^{13}$ unconstrained (OpenTuner) vs. $10^6$ constrained (ATF).

  - XgemmDirect has to rely on its default tuning parameter values → chosen to yield a good performance on average on various devices and for different input sizes.
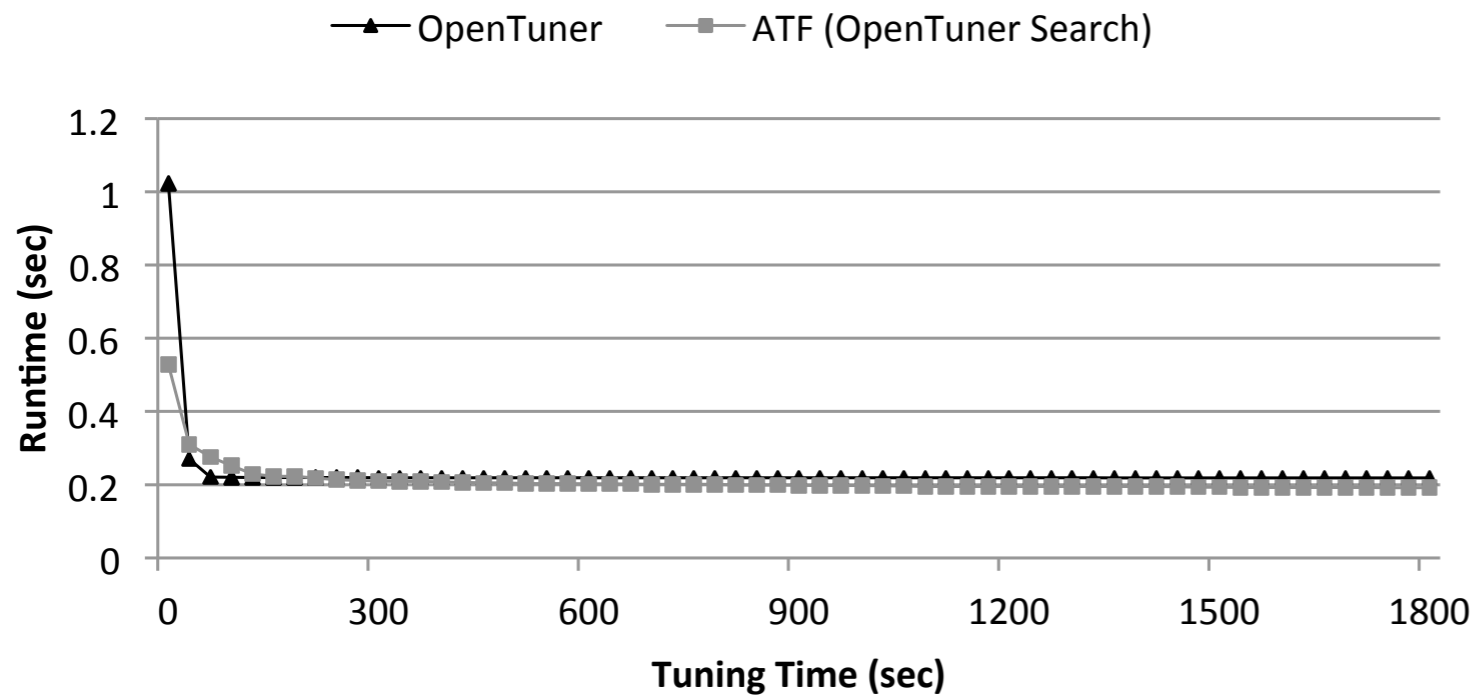
# Summary

- Auto-tuning simplifies optimizing programs by automatically choosing suitable values of tuning parameters.

- ATF provides four advantages over the state-of-the-art auto-tuning approaches:

  1. ATF is **generic** regarding the programming language, application domain, tuning objective, and search technique.

  2. ATF allows **dependencies between tuning parameters**, thus enabling to auto-tune a broader class of applications.

  3. ATF allows **significantly larger tuning parameter ranges** and thus do not require artificially limiting parameters' ranges $\Rightarrow$ better performance.

  4. ATF is **arguably simpler** to use, thus making auto-tuning appealing to common application developers.

- ATF significantly accelerates the performance of GEMM on practically-relevant input sizes as used in deep learning.

# Experimental Results

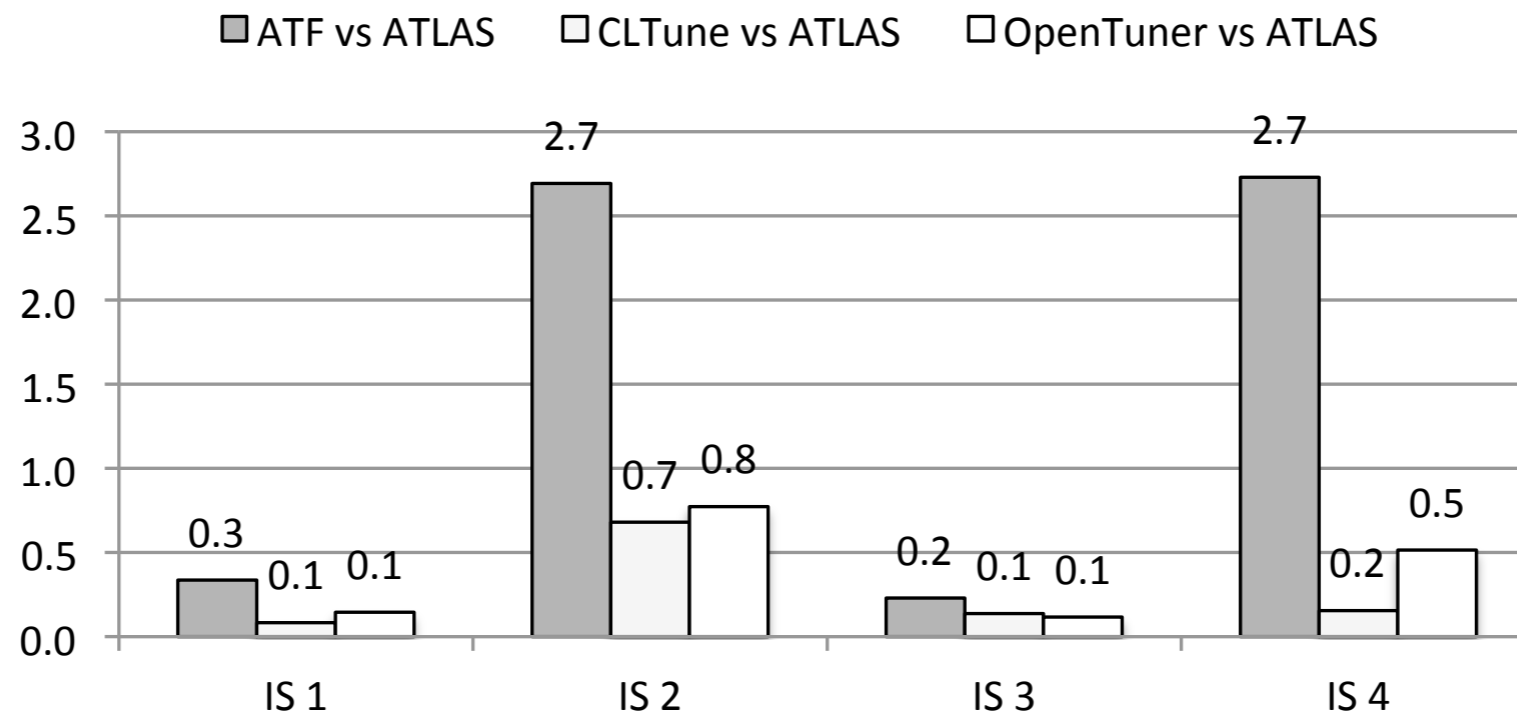ATF provides results of competitive quality as compared to CLTune and OpenTuner for their target application classes:



## 2D Convolution (CLTune)



## GCC Flags (OpenTuner)

# Appendix

## ATF vs ATLAS:



- In constrast to CLTune and OpenTuner, ATF is able to auto-tune XgemmDirect to better performance than ATLAS for the input sizes IS 2 and IS 4 --- a speedup of 2.7x in both cases.

- This again due to ATF's support for large parameter ranges (CLTune) and ATF's parameter constraints (OpenTuner).

- ATF was not able to auto-tune XgemmDirect to better performance than ATLAS for the input sizes IS 1 and IS 3.

- We argue that this is due to the implementation of XgemmDirect which is not optimized for the special matrix input matrices whose number of rows/columns are 1 ($\rightarrow$ IS 1: 20×1 * 1×576; IS 3: 50×1 * 1×64).

# Appendix

Search space generation:

```
for( val_1 : tp_1.range )
  if( constraint_1( val_1 ) == true )
    .
     .
      .
       for( val_n : tp_n.range )
         if( constraint_n( val_n ) == true )
         {
           search_space.add( val_1, ... , val_n );
         }
```

ATF

```
for( val_1 : tp_1.range )
  .
   .
    .
     for( val_n : tp_n.range )
     {
       for( c : constraints )
         if( c( val_1, ... , val_n )
           search_space.add( val_1, ... , val_n );
     }
```

CLTune

# Appendix

ATF's "*search_technique*" interface:

```
class search_technique
{
  void          initialize( search_space sp );
  void          finalize();
  configuration get_next_config();
  void          report_cost( size_t cost );
}
```