# Autotuning Search Space for Loop Transformations

6[th] Workshop on the LLVM Compiler Infrastructure in HPC

Michael Kruse, Xingfu Wu, Hal Finkel

Argonne National Laboratory

2020-11-12

# Outline

# Outline

# Partial Unrolling

```
#pragma unroll 4
for (int i = 0; i < n; i += 1)
  Stmt(i);



if (n > 0) {
  for (int i = 0; i+3 < n; i += 4) {
    Stmt(i);
    Stmt(i + 1);
    Stmt(i + 2);
    Stmt(i + 3);
  }
  switch (n % 4) {
  case 3:
    Stmt(n - 3);
  case 2:
    Stmt(n - 2);
  case 1:
    Stmt(n - 1);
  }
}
```

# Compiler-Supported Pragmas

**Clang**
```
#pragma unroll
#pragma clang loop unroll(enable)
#pragma unroll_and_jam
#pragma clang loop distribute(enable)
#pragma clang loop vectorize(enable)
#pragma clang loop interleave(enable)
#pragma clang loop pipeline(disable)
```

**gcc**
```
#pragma GCC unroll
#pragma GCC ivdep
```

**msvc**
```
#pragma loop(hint_parallel(0))
#pragma loop(no_vector)
#pragma loop(ivdep)
```

**Cray**
```
#pragma _CRI unroll
#pragma _CRI fusion
#pragma _CRI nofission
#pragma _CRI blockingsize
#pragma _CRI interchange
#pragma _CRI collapse
```

**OpenMP**
```
#pragma omp simd
#pragma omp for
#pragma omp target
```

**PGI**
```
#pragma concur
#pragma vector
#pragma ivdep
#pragma nodepchk
```

**xlc**
```
#pragma unrollandfuse
#pragma stream_unroll
#pragma block_loop
#pragma loopid
```

**SGI/Open64**
```
#pragma fuse
#pragma fission
#pragma blocking size
#pragma altcode
#pragma noinvarif
#pragma mem prefetch
#pragma interchange
#pragma ivdep
```

**OpenACC**
```
#pragma acc kernels
```

**icc**
```
#pragma parallel
#pragma offload
#pragma unroll_and_jam
#pragma nofusion
#pragma distribute_point
#pragma simd
#pragma vector
#pragma swp
#pragma ivdep
#pragma loop_count(n)
```

**Oracle Developer Studio**
```
#pragma pipeloop
#pragma nomemorydepend
```

**HP**
```
#pragma UNROLL_FACTOR
#pragma IF_CONVERT
#pragma IVDEP
#pragma NODEPCHK
```

# Loop Transformation Zoo

### Tiling

```
#pragma clang transform tile sizes(4,4)
for (int i = 0; i < m; i += 1)
  for (int j = 0; j < n; j += 1)
    Body(i,j);
```

```
for (int i1 = 0; i1 < m; i1 += 4)
  for (int j1 = 0; j1 < n; j1 += 4)
    for (int i2 = i1; i2 < min(i1+4,m); i2 += 1)
      for (int j2 = i1; j2 < min(j1+4,n); j2 += 1)
        Body(i2,j2);
```

### Unrolling

```
#pragma clang transform unroll partial(4)
for (int i = 0; i < n; i += 1)
  Body(i);
```

```
for (int i = 0; i < n; i1 += 4) {
  Body(i);
  Body(i+1);
  Body(i+2);
  Body(i+3);
}
```

### Fusion

```
#pragma clang transform fuse
for (int i = 0; i < n; i+=1)
  BodyA(i);
for (int i = 0; i < n; i+=1)
  BodyB(i);
```

```
for (int i = 0; i < n; i+=1) {
  BodyA(i);
  BodyB(i);
}
```

### Space-Filling Curves

```
#pragma clang transform spacefill curve(hilbert)
for (int i = 0; i < m; i += 1)
  for (int j = 0; j < n; j += 1)
    Body(i,j);
```

```
for (int idx = 0; idx < m*n; idx += 1) {
  tie(i,j) = hilbert2d_from_index(idx,m,n);
  Body(i,j);
}
```

### Interchange

```
#pragma clang transform interchange
for (int i = 0; i < m; i+=1)
  for (int j = 0; j < n; j+=1)
    Body(i,j);
```

```
for (int j = 0; j < n; j+=1)
  for (int i = 0; i < m; i+=1)
    Body(i,j);
```

### Reversal

```
#pragma clang transform reverse
for (int i = 0; i < n; i+=1)
  Body(i);
```

```
for (int i = n-1; i >= 0; i-=1)
  Body(i);
```

# Composition of Transformations

```
#pragma unroll 2
#pragma reverse
for (int i = 0; i < 128; i+=1)
  Stmt(i);
```

⬇

```
#pragma unroll 2
for (int i = 127; i >= 0; i-=1)
  Stmt(i);
```

⬇

```
for (int i = 127; i >= 0; i-=1) {
  Stmt(i);
  Stmt(i-1);
}
```

Stmt(127); Stmt(126); Stmt(125); Stmt(124); ...

```
#pragma reverse
#pragma unroll 2
for (int i = 0; i < 128; i+=1)
  Stmt(i);
```

⬇

```
#pragma reverse
for (int i = 0; i < 128; i+=2) {
  Stmt(i);
  Stmt(i+1);
}
```

⬇

```
for (int i = 126; i >= 0; i-=2) {
  Stmt(i);
  Stmt(i+1);
}
```

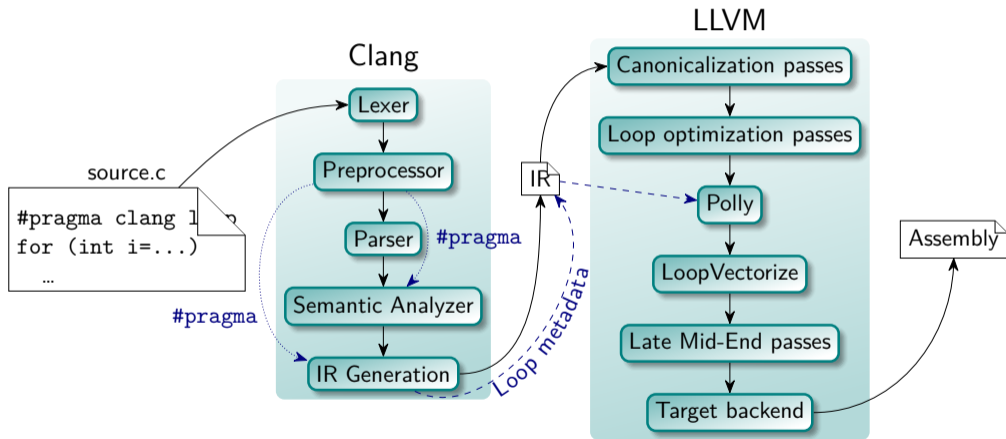Stmt(126); Stmt(127); Stmt(124); Stmt(125); ...

# Advanced Loop Directives Prototype

- *#pragma clang loop id(...)*
  Name a loop to disambiguate transformation order
- *#pragma clang loop reverse*
  Loop Reversal
- *#pragma clang loop tile sizes(...)*
  Tiling
- *#pragma clang loop interchange permutation(...)*
  Loop interchange/permutation
- *#pragma clang loop unrolling factor(...)*
  Loop Unrolling
- *#pragma clang loop unrollingandjam factor(...)*
  Loop Unroll-and-Jam
- *#pragma clang loop parallelize_thread*
  Thread-Parallelization
- *#pragma clang loop pack array(....)*
  Array Packing

# Prototype Implementation

## Implementations

- *#pragma clang loop(id) <transformation> <clauses...>*
  (Prototype implementation, LLVM-in-HPC'18)
  `https://github.com/sollve/llvm-project/tree/pragma-clang-loop`

- *#pragma clang transform <transformation> on(id) <clauses...>*
  (Composable transformations in Clang)
  `https://reviews.llvm.org/D69088`

- *#pragma clang omp <construct> <clauses...>*
  (OpenMP 5.1 implementation for Clang)
  `https://reviews.llvm.org/D76342`

## PolyBench/C dgemm

```c
void matmul(int M, int N, int K, double alpha, double beta,
            double C[const restrict static M][N],
            double A[const restrict static M][K],
            double B[const restrict static K][N]) {

  for (int i = 0; i < M; i++)
    for (int j = 0; j < N; j++) {
      C[i][j] *= beta;
      for (int k = 0; k < K; k++)
        C[i][j] += alpha * A[i][k] * B[k][j];
    }

}
```

## PolyBench/C dgemm

```
void matmul(int M, int N, int K, double alpha, double beta,
            double C[const restrict static M][N],
            double A[const restrict static M][K],
            double B[const restrict static K][N]) {

  for (int i = 0; i < M; i++)
    for (int j = 0; j < N; j++)
      C[i][j] *= beta;

  for (int i = 0; i < M; i++)
    for (int j = 0; j < N; j++)
      for (int k = 0; k < K; k++)
        C[i][j] += alpha * A[i][k] * B[k][j];

}
```

# PolyBench/C dgemm
Using Inner Loop Vectorization

```
#pragma clang loop(j2) pack array(A)
#pragma clang loop(i1) pack array(B)
#pragma clang loop(i1,j1,k1,i2,j2) interchange permutation(j1,k1,i1,j2,i2)
#pragma clang loop(i,j,k) tile sizes(96,2048,256) \
                          floor_ids(i1,j1,k1) tile_ids(i2,j2,k2)

#pragma clang loop id(i)
for (int i = 0; i < M; i += 1)
  #pragma clang loop id(j)
  for (int j = 0; j < N; j += 1)
    #pragma clang loop id(k)
    for (int k = 0; k < K; k += 1)
      C[i][j] += alpha * A[i][k] * B[k][j];
```

# Polybench dgemm
## After Transformation (Sanitized)

```
double Packed_B[256][2048];
double Packed_A[96][256];
if (runtime check) {
  if (M >= 1)
    for (int c0 = 0; c0 <= floord(N - 1, 2048); c0 += 1)     // Loop j1
      for (int c1 = 0; c1 <= floord(K - 1, 256); c1 += 1) { // Loop k1

        // Copy-in: B -> Packed_B
        for (int c4 = 0; c4 <= min(2047, N - 2048 * c0 - 1); c4 += 1)
          for (int c5 = 0; c5 <= min(255, K - 256 * c1 - 1); c5 += 1)
            Packed_B[c4][c5] = B[256 * c1 + c5][2048 * c0 + c4];

        for (int c2 = 0; c2 <= floord(M - 1, 96); c2 += 1) { // Loop i1

          // Copy-in: A -> Packed_A
          for (int c6 = 0; c6 <= min(95, M - 96 * c2 - 1); c6 += 1)
            for (int c7 = 0; c7 <= min(255, K - 256 * c1 - 1); c7 += 1)
              Packed_A[c6][c7] = A[96 * c2 + c6][256 * c1 + c7];

          for (int c3 = 0; c3 <= min(2047, N - 2048 * c0 - 1); c3 += 1)   // Loop j2
            for (int c4 = 0; c4 <= min(95, M - 96 * c2 - 1); c4 += 1)     // Loop i2
              for (int c5 = 0; c5 <= min(255, K - 256 * c1 - 1); c5 += 1) // Loop k2
                C[96 * c2 + c4][2048 * c0 + c3] += Packed_A[c4][c5] * Packed_B[c3][c5];
        }
      }
} else {
  /* original code */
}
```

## Polybench dgemm



Double precsion FP operations per time unit in percentage of peak

## Motivation

- User-directives requires the user to know what to do …
- … why not find the transformations automatically?
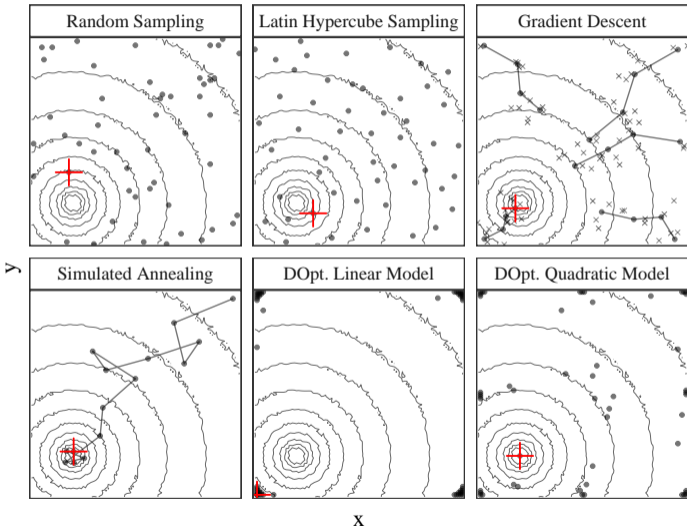
# Outline

# Vector Space Combinatorial Optimization

Predictive Modeling / Empirical



P. Bruel et. al.
Autotuning under Tight Budget Constraints:
A Transparent Design of Experiments Approach
(CCGRID'19)

## Components

- **ytopt**: Autotuning framework, supports multiple machine learning strategies
  https://github.com/ytopt-team/ytopt
- **Plopper**: Source rewrite engine
  https://github.com/ytopt-team/autotune
- **pragma-clang-loop**: Loop transformation implementation
  https://github.com/sollve/llvm-project/tree/pragma-clang-loop

# Plopper
syr2k

```
#P0
#P1
#P2
#pragma clang loop(i,j,k) tile sizes(#P3,#P4,#P5) \
        floor_ids(i1,j1,k1) tile_ids(i2,j2,k2)

#pragma clang loop id(i)
for (i = 0; i < _PB_N; i++)
    #pragma clang loop id(j)
    for (j = 0; j < _PB_M; j++)
     #pragma clang loop id(k)
        for (k = 0; k <= i; k++)
          C[i][k] += A[k][j]*alpha*B[i][j] + B[k][j]*alpha*A[i][j];
```

# ytopt Search Space Definition

```
P0=CSH.CategoricalHyperparameter(name='P0', default_value=' ',
    choices=["#pragma clang loop(j2) pack array(A) allocate(malloc)", " "])
P1=CSH.CategoricalHyperparameter(name='P1', default_value=' ',
    choices=["#pragma clang loop(i1) pack array(B) allocate(malloc)", " "])
P2=CSH.CategoricalHyperparameter(name='P2', default_value=' ',
    choices=["#pragma clang loop(i1,j1,k1,i2,j2) interchange permutation(j1,k1,i1,j2,i2)", " "])
P3=CSH.OrdinalHyperparameter(name='P3', default_value='96',
    sequence=['4','8','16','20','32','50','64','80','96','100','128'])
P4=CSH.OrdinalHyperparameter(name='P4', default_value='2048',
     sequence=['4','8','16','20','32','50','64','80','100','128','2048'])
P5=CSH.OrdinalHyperparameter(name='P5',  default_value='256',
    sequence=['4','8','16','20','32','50','64','80','100','128','256'])
cs.add_hyperparameters([P0, P1, P2, p3, P4, P5])
cond1 = CS.InCondition(P1, P0, ['#pragma clang loop(j2) pack array(A) allocate(malloc)'])
```

# Syr2k Tuning
Random Forests



Autotuning syr2k Using Random Forests

# Syr2k Tuning Result

```
#pragma clang loop(j2) pack array(A) allocate(malloc)
#pragma clang loop(i1) pack array(B) allocate(malloc)
#pragma clang loop(i1,j1,k1,i2,j2) interchange permutation(j1,k1,i1,j2,i2)
#pragma clang loop tile sizes(128,128,100) \
        floor_ids(i1,j1,k1) tile_ids(i2,j2,k2)
for (int i = 0; i < _PB_N; i++)
    for (int j = 0; j < _PB_M; j++)
        for (int k = 0; k <= i; k++)
            C[i][k] += A[k][j]*alpha*B[i][j] + B[k][j]*alpha*A[i][j];
```

# Vector-Space Autotuning

## Strengths

- Machine learning over tuples of numbers well established
- Exploit contiguous relationships

## Weaknesses

- User has to define loop transformations
- Search space is necessarily incomplete
  (Number of transformations is potentially infinite)
- Depending on previous choices, some parameters are unused
- Some parameters cause abrupt changes
  (e.g. no predictable relationship between loop permutations)

# Related Work

## CHiLL (TVM, …)

- Transformations in hand-written *recipes*
- Tune placeholders in recipe (e.g. tile sizes)

## LIFT (SPIRAL,…)

- Based on rewriting engine
- Using ATF (Auto-Tuning Framework)
- Parameter constraints to prune search space

## LeTSeE

- Encode schedule function coefficients
- Correctness constrains encoded as polytope boundaries

## PATUS (CLTuner,…)

- Specific to categories of kernels (stencils)
- Assume trivially parallel kernels
- Tune fixed set of parameters (tile size, workgroup size, …)

## HalideTuner

- Using OpenTuner
- Encode transformation sequence in constant-size vector
- Verifier prunes configuration that do no encode a valid transformation sequence

# Outline

## Successive Optimization

```
for (int i = 0; i < M; i += 1)
  for (int j = 0; j < N; j += 1)
    for (int k = 0; k < K; k += 1)
      C[i][j] += A[i][k] * B[k][j];
```
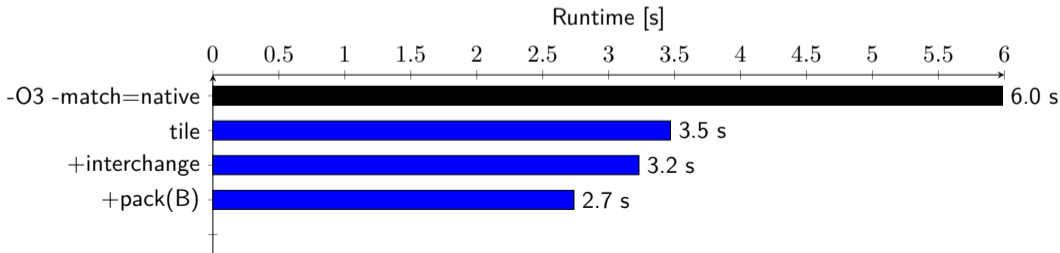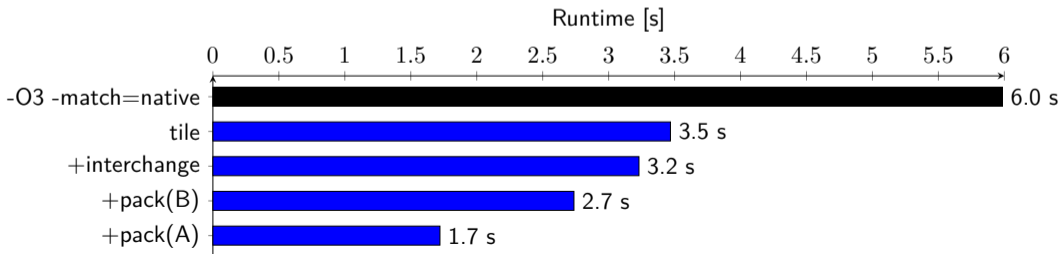
Runtime [s]

| 0 | 0.5 | 1 | 1.5 | 2 | 2.5 | 3 | 3.5 | 4 | 4.5 | 5 | 5.5 | 6 |

-O3 -match=native  ████████████████████████████████████████ 6.0 s

## Successive Optimization

```
#pragma clang loop tile sizes(96,2048,256) floor_ids(i1,j1,k1) tile_ids(i2,j2,k2)
for (int i = 0; i < M; i += 1)
  for (int j = 0; j < N; j += 1)
    for (int k = 0; k < K; k += 1)
      C[i][j] += A[i][k] * B[k][j];
```
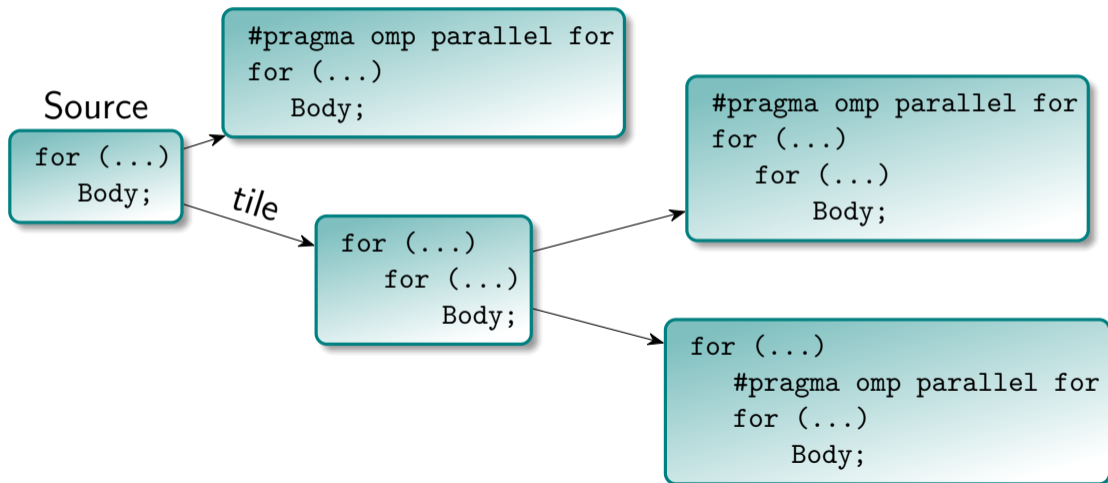
## Successive Optimization

```
#pragma clang loop(i1,j1,k1,i2,j2) interchange permutation(j1,k1,i1,j2,i2)
#pragma clang loop tile sizes(96,2048,256) floor_ids(i1,j1,k1) tile_ids(i2,j2,k2)
for (int i = 0; i < M; i += 1)
  for (int j = 0; j < N; j += 1)
    for (int k = 0; k < K; k += 1)
      C[i][j] += A[i][k] * B[k][j];
```
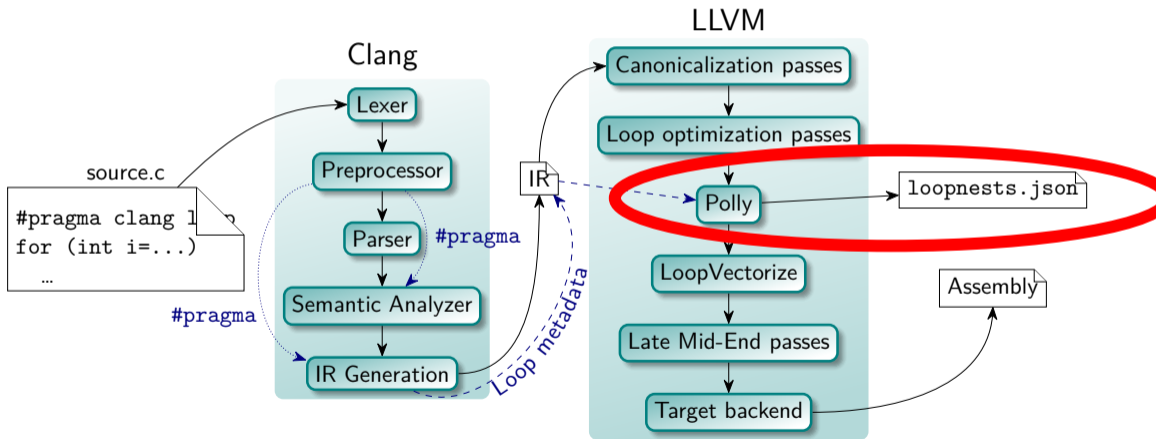
## Successive Optimization

```
#pragma clang loop(i1) pack array(B)
#pragma clang loop(i1,j1,k1,i2,j2) interchange permutation(j1,k1,i1,j2,i2)
#pragma clang loop tile sizes(96,2048,256) floor_ids(i1,j1,k1) tile_ids(i2,j2,k2)
for (int i = 0; i < M; i += 1)
  for (int j = 0; j < N; j += 1)
    for (int k = 0; k < K; k += 1)
      C[i][j] += A[i][k] * B[k][j];
```

## Successive Optimization

```
#pragma clang loop(j2) pack array(A)
#pragma clang loop(i1) pack array(B)
#pragma clang loop(i1,j1,k1,i2,j2) interchange permutation(j1,k1,i1,j2,i2)
#pragma clang loop tile sizes(96,2048,256) floor_ids(i1,j1,k1) tile_ids(i2,j2,k2)
for (int i = 0; i < M; i += 1)
  for (int j = 0; j < N; j += 1)
    for (int k = 0; k < K; k += 1)
      C[i][j] += A[i][k] * B[k][j];
```

# Transformation Composition Possibilities



Source

```
for (...)
    Body;
```

```
#pragma omp parallel for
for (...)
    Body;
```

*tile*

```
for (...)
    for (...)
        Body;
```

```
#pragma omp parallel for
for (...)
    for (...)
        Body;
```

```
for (...)
    #pragma omp parallel for
    for (...)
        Body;
```

# Prototype Implementation

# Implementation: Start

1. Compile
   - Output: Loop nest structure in `loopnests.json`
2. Run
   - Measure baseline wall time
3. Generate search tree ("*mctree*")

## Proof-of-Concept Transformations

- Tiling
  - Tile sizes: 4, 16, 64, 256, 1024
- Interchange
- Thread parallelization (`schedule(static)`)

# Implementation: Iterate

**1** Select configuration to run
  - According search strategy

**2** Insert transformation directives

**3** Compile

**4** Run
  - Measure baseline wall time

**5** Rinse, Repeat

## Proof-of-Concept Search Strategy

- Select child of fastest-yet

# Tree Search Space Autotuning

# Tree Search Space Autotuning

- PolyBench/C 4.2 kernels: dgemm, syr2k, covariance
    - `EXTRALARGE_DATASET` (largest predefined problem size)
    - Manually distributed
    - syr2k and covariance are non-rectangular

- 2x Intel Xeon Platinum 8180M
    - 28 cores each, 112 logical threads
    - 376 GiB RAM, 38.5 MiB L3, 1 MiB L2, 32 KiB L1d

- 6h tuning time
    - Program execution time limit: 120 secs.

# dgemm Tuning
Interchange/Tile/Parallel



Best configuration:

# dgemm Tuning
Interchange/Tile/Parallel



Best configuration:  `#pragma clang loop(i) parallelize_thread`
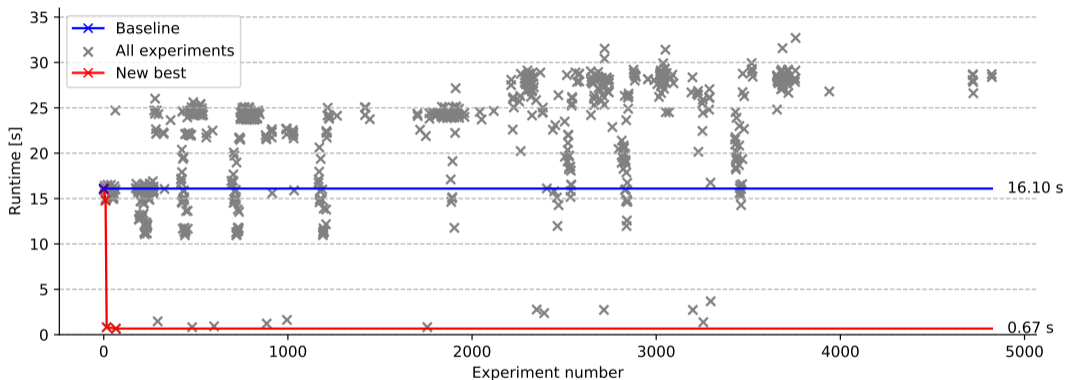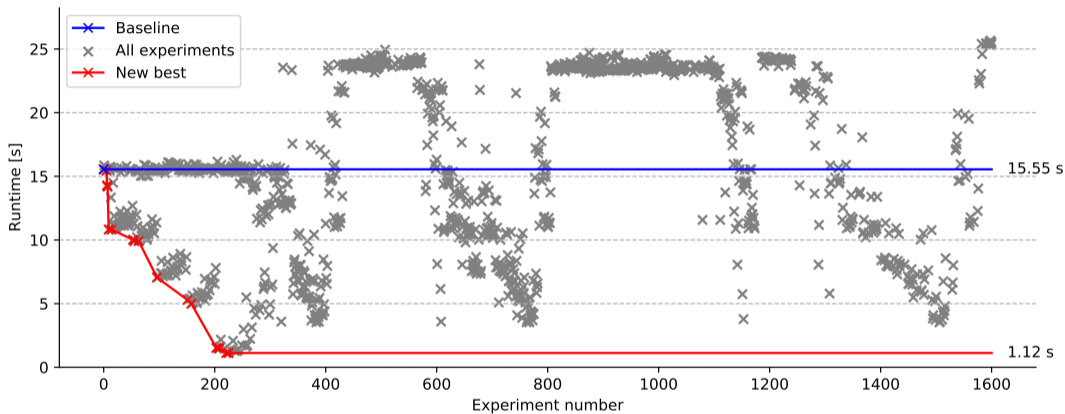
# dgemm Tuning
Interchange/Tile



Best configuration:

```
#pragma clang loop(j,k,i) tile sizes(1024,64,16)
#pragma clang loop(i,j,k) interchange permutation(j,k,i)
```

# syr2k Tuning
Interchange/Tile/Parallel



Best configuration: *#pragma clang loop(i) parallelize_thread*

# syr2k Tuning
Interchange/Tile



Best configuration:  *#pragma clang loop(i,j,k) tile sizes(64,256,4)*

# covariance Tuning Result

Interchange/Tile/Parallel



Best configuration:   *#pragma clang loop(i) parallelize_thread*
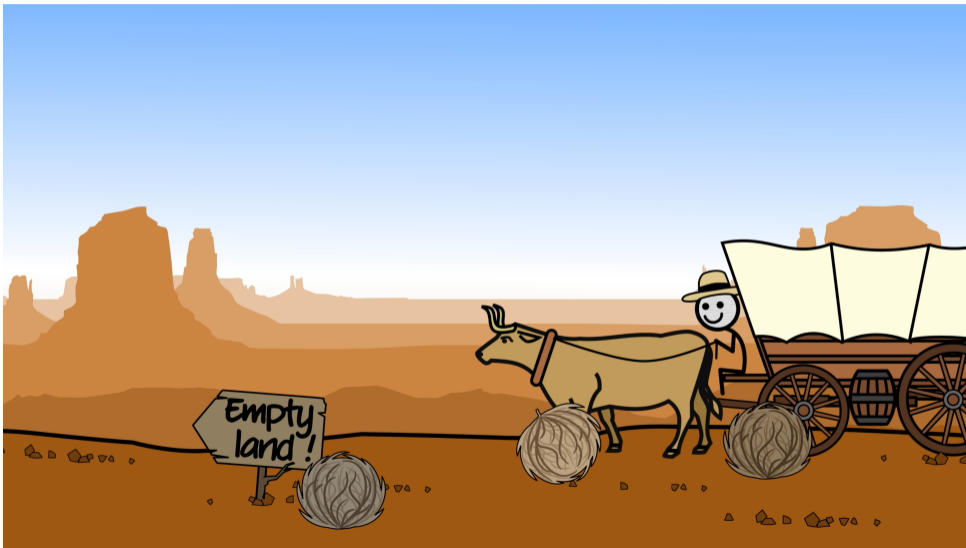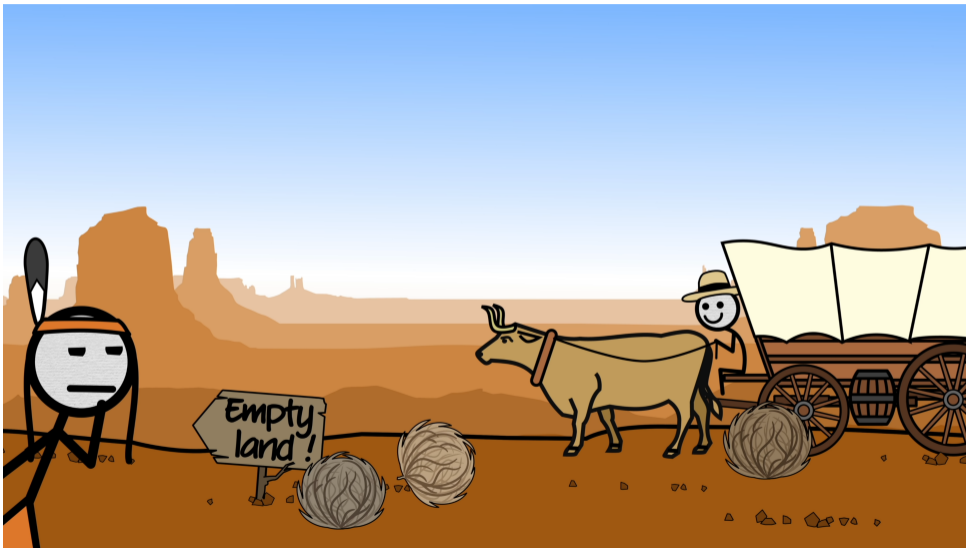
# covariance Tuning Result

Interchange/Tile



Best configuration:
$$\#pragma\ clang\ loop(j,k,i)\ tile\ sizes(1024,64,16)$$
$$\#pragma\ clang\ loop(i,j,k)\ interchange\ permutation(j,k,i)$$

# Similar Work



youtu.be/hsWr_JWTZss
CGP Grey: The Trouble With Tumbleweed

# Similar Work



youtu.be/hsWr_JWTZss
CGP Grey: The Trouble With Tumbleweed

# Similar Work

## Halide (ACM TOG, July 2019)

`https://dl.acm.org/doi/abs/10.1145/3306346.3322967`

- Use transformation tree space (organized in stages)
- Beam Search
- Magnitudes faster than HalideTune
- Twice as fast as Halide's latest more restricted default autotuner.

## ProTuner (arXiv, May 2020)

`https://arxiv.org/abs/2005.13685`

- Use same transformation tree space
- Monte-Carlo tree search
- Up to 3.6 times faster executable that with Beam Search

# Outline

# Future Work / Summary

## Summary

- Two search space shapes
  - Vector space
  - Tree space
- Implement loop transformation config space generator (*mctree*)
- Implement proof-of-concept autotuner
  - Greedy search space exploitation
- Apply to PolyBench kernels
  - gemm, syr2k, covariance

## Future Work

- Implement Monte-Carlo tree search
  - Do not get stuck in local optima
- Prune search space
  - Interchange twice is redundant
  - Multiple paths to same configuration ($\Rightarrow$ DAG)
- Combine approaches
  - MCTS for transformation composition
  - ytopt for transformation parameters
- Measure/tune loop nests individually
  - Avoid combinatorial explosion

## Acknowledgments