# Data Parallel Programming in Java Using the Vector API

Adam Pocock, Oracle
Paul Sandoz, Oracle

JavaOne™

# Overview

Motivation

Use cases

Concepts

Examples

Guidance

Status

# Motivation

- Most modern CPUs have

  — Special registers called vector registers that can hold multiple data

  — Many vector instructions that operate on those registers

- A **S**ingle **I**nstruction can operate on **M**ultiple **D**ata (**SIMD**)

  — We can use these instructions to implement highly performant data parallel algorithms

- How can we leverage all those instructions using Java?

# Motivation

- OpenJDK's HotSpot runtime compiler can leverage some of these instructions

  - Some Java loops can be transformed into data-parallel loops when generating machine code. This is process is called **auto-vectorization**

  - HotSpot is very good at this, but the technique is limited

- How can we overcome the limitations of HotSpot's auto-vectorization?

  - To reliably express a broad set of data parallel algorithms?

# Vector API

- Enables a developer to reliably express cross-platform data parallel algorithms leveraging a significant set of vector instructions

  — A **W**hat **Y**ou **S**ee is **W**hat **Y**ou **G**et (**WYSIWYG**) API, almost

- HotSpot reliably compiles code written using the Vector API to vector instructions

  — On Intel, ARM, Power CPUs, and soon RISC-V CPUs

  — Hiding the various differences in CPU architectures

# Use cases

- Daniel Lemire's blog and publications present some excellent use cases and highly performant algorithms

  - Lemire's JavaFastPFOR project recently integrated an implementation using the Vector API

- In OpenJDK

  - String encoding, BASE64 encoding, array/string comparison, array hash codes, sorting, cryptography

  - Some are currently supported with explicit assembler code

# Example use cases

- Filtering

- Machine learning

- (Lanewise arc-tangent)

# Concepts

- A `Vector` is a sequence of elements. Each element is held in a lane
  - Many operations can be performed on vectors

- A `VectorMask` is a sequence of booleans
  - Masks can control if an operation is applied to a lane

- A `VectorShuffle` is a sequence of lane indexes
  - Shuffles can rearrange the elements of a vector

- All sequences are **logical**, fixed in size, ordered, and immutable

# Concepts

- A `VectorSpecies` manages vectors, masks, and shuffles
  - Species can instantiate vectors with a specific element type and shape that together determines the length of the vector (or number of lanes)

- A `VectorShape` selects a particular implementation of vectors, masks and shuffles
  - A shape of 256 bits selects vectors whose elements can be stored in a vector register of size 256 bits or greater
  - A species with a 256 bits shape and an element type of `float` will instantiate vectors with 8 lanes

# Concepts

- Lanewise operations

  — Unary, binary, ternary, unary test, compare (binary test)

- Cross-lane operations

  — Reduction, rearranging, slicing, compressing/expanding

- Reinterpret operations

  — Conversion between vectors of different element types and shapes

# Compute the sum of all positive elements of an array

```
1   @Benchmark
2   public float scalarf() {
3       float sum = 0.0f;
4       for (int i = 0; i < fvalues.length; i++) {
5           float v = fvalues[i];
6           if (v > 0) sum += v;
7       }
8
9       return sum;
10  }
```

# Data parallel code using the Vector API

```java
1   static final VectorSpecies<Float> F_S = FloatVector.SPECIES_PREFERRED;
2
3   @Benchmark
4   public double vectorf() {
5       FloatVector vsum = FloatVector.broadcast(F_S, 0.0f);
6       int i = 0;
7       for (; i < F_S.loopBound(fvalues.length); i += F_S.length()) {
8           FloatVector v = FloatVector.fromArray(F_S, fvalues, i);
9           VectorMask<Float> m = v.compare(VectorOperators.GT, 0.0f);
10          vsum = vsum.add(v, m);
11      }
12      float sum = vsum.reduceLanes(VectorOperators.ADD);
13      for (; i < fvalues.length; i++) {
14          float v = fvalues[i];
15          if (v > 0) sum += v;
16      }
17      return sum;
18  }
```

# Data parallel code using the Vector API

```
species: shape = 256 bits, element type = float
       ∴ lanes = 256 / 32 = 8

// VectorMask<Float> m = v.compare(VectorOperators.GT, 0.0f);
v[1.0, -2.0, 3.0, -4.0, 5.0, -6.0, 7.0, -8.0]
    >     >    >     >    >     >    >     >
 [0.0,  0.0, 0.0,  0.0, 0.0,  0.0, 0.0,  0.0])
    =     =    =     =    =     =    =     =
m[  T,    F,   T,    F,   T,    F,   T,    F]

// vsum = vsum.add(v, m);
vsum[1.0,   1.0,   1.0,   1.0, 1.0,   1.0, 1.0,   1.0]
   m[  T,     F,     T,     F,   T,     F,   T,     F]
       +             +           +           +
   v[1.0, -2.0,   3.0, -4.0, 5.0, -6.0, 7.0, -8.0]
       =             =           =           =
vsum[2.0,   1.0,   4.0,   1.0, 6.0,   1.0, 8.0,   1.0]
```

# Generated code is good

```
// 2.6 GHz 6-Core Intel Core i7
// Produced using JMH's dtraceasm profiler
// -XX:-TieredCompilation -XX:LoopUnrollLimit=0

  0.14%   ↗      0x000000011e9726b0:    vmovdqu 0x10(%rsi,%rdi,4),%ymm1
 10.73%          0x000000011e9726b6:    vaddps %ymm1,%ymm0,%ymm3
 35.68%          0x000000011e9726ba:    vcmpgt_oqps %ymm2,%ymm1,%ymm1
  1.00%          0x000000011e9726bf:    vblendvps %ymm1,%ymm3,%ymm0,%ymm0
 42.19%          0x000000011e9726c5:    add     $0x8,%edi
  2.14%          0x000000011e9726c8:    cmp     %r10d,%edi
                 0x000000011e9726cb:    jl      0x000000011e9726b0
```

# Comparing source with generated code (WYSIWYG)

```
1   for (;
2        // cmp      %r10d,%edi
3        i < F_S.loopBound(dvalues.length);
4        // add      $0x8,%edi
5        i += F_S.length()) {
6
7        // vmovdqu
8        FloatVector v = FloatVector.fromArray(F_S, fvalues, i);
9        // vcmpgt_oqps
10       VectorMask<Float> m = v.compare(VectorOperators.GT, 0.0f);
11       // vaddps
12       // vblendvps
13       vsum = vsum.add(v, m);
14   }
```

# Performance is good (lower is better)

```
// 2.6 GHz 6-Core Intel Core i7
// -XX:-TieredCompilation
Benchmark   (size)   Mode   Cnt    Score     Error   Units
scalarf       1024   avgt     5  982.786 ± 44.711   ns/op
vectorf       1024   avgt     5  169.881 ±  4.465   ns/op
```

# Example: filtering

- Imagine a table arranged as an in-memory columnar data structure

```
position,  value,    x,    y,    z
      p0      v0   ...,  ...,  ...
      p1      v1
      p2      v2
      p3      v3
      p4      v4
      p5      v5
      p6      v6
      p7      v7
      ...     ...
```

- The first column represents the position or identifier of rows in the table

# Filtering

- How can we filter the rows whose `values` are within some upper and lower bound?

e.g.,

```
positions[p0, p1, p2, p3, p4, p5, p6, p7]
   values[-1,  0,  5,  4,  3,  2, 10, 12]
          0 <= v <= 4
 filtered[p1, p3, p4, p5]
```

# Sequential code

```
1   public boolean testInt(int value) {
2       return value >= lower && value <= upper;
3   }
4
5   @Benchmark
6   public void scalari() {
7       int matchCount = 0;
8       for (int i = 0; i < values.length; i++) {
9           if (testInt(values[i])) {
10              filtered[matchCount++] = positions[i];
11          }
12      }
13  }
```

# Data parallel code using the Vector API

- The Vector API in JDK 19 has a new cross-lane operation, `compress` (and its inverse `expand`)

  – This operation is optimized on supporting hardware e.g. Intel CPUs supporting the AVX-512 instruction set

- We can convert the sequential code to data parallel code using this new operation

# Data parallel code using the Vector API

```java
1  public VectorMask<Integer> testIntVector(IntVector values) {
2      return values.compare(GE, lower).and(values.compare(LE, upper));
3  }
4
5  @Benchmark
6  public void vectori() {
7      int matchCount = 0, i = 0;
8      for (; i < I_S.loopBound(values.length); i += I_S.length()) {
9          var v = IntVector.fromArray(I_S, values, i);
10         VectorMask<Integer> match = testIntVector(v);
11
12         var pV = IntVector.fromArray(I_S, positions, i);
13         var fV = pV.compress(match);
14
15         fV.intoArray(filtered, matchCount);
16         matchCount += match.trueCount();
17     }
18     ...
```

# Data parallel code using the Vector API

```
// VectorMask<Integer> match = testIntVector(v);
values[-1,  0,  5,  4,  3,  2, 10, 12]
       0 <= v <= 4
 match[ F,  T,  F,  T,  T,  T,  F,  F]

// var fV = pV.compress(match);
positions[p0, p1, p2, p3, p4, p5, p6, p7]
    match[ F,  T,  F,  T,  T,  T,  F,  F]
          compress
 filtered[p1, p3, p4, p5,  0,  0,  0,  0]
```

# Performance is good (higher is better)

```
// On Intel AVX-512 system
Benchmark (selectivity)    Mode    Cnt          Score            Error  Units
scalari                128  thrpt      3   1294956.909 ±     1463.924  ops/s
vectori                128  thrpt      3   7574101.344 ±    30222.336  ops/s
```

# Sorting

- We can use `compress` to optimize Quicksort
  - Quicksort recursively partitioning elements around a pivot point
  - See [blog](#) and [paper](#) by Google engineers
- Applicable to Java's sorting of primitive arrays

# Example: machine learning

- Machine Learning is a heavy consumer of FLOPS, which the Vector API makes much more accessible in Java

- There are a few basic operations:
  - Matrix multiply
  - Convolution
  - Non-linear element-wise functions like max(0,x), sigmoid, or tanh

- We'll look at the 7x7 convolution operation used as the first step in most image recognition or object detection systems

# Convolution

- The convolution operation multiplies each "patch" of an image with a learned weight matrix - the "convolution matrix"

- We move the learned matrix across the image one pixel at a time, multiplying each element of the convolution with the covered image pixel and summing them to produce an output pixel

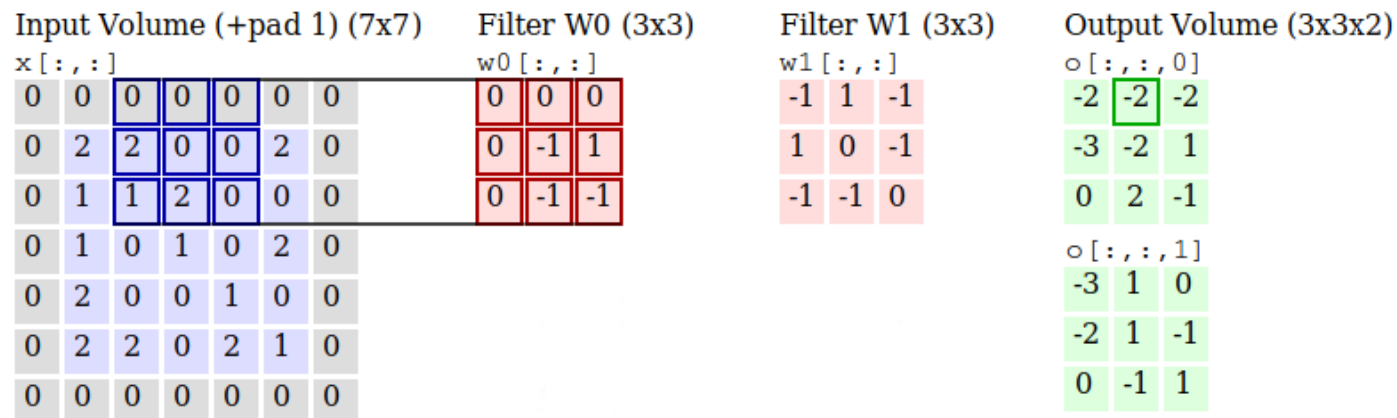- It is essentially a per pixel dot product, with some fancy indexing



| Input Volume (+pad 1) (7x7) | Filter W0 (3x3) | Filter W1 (3x3) | Output Volume (3x3x2) |
|---|---|---|---|

Figure adapted from Stanford's CS231n course

# Scalar convolution

```java
1   // Compute convolution
2   for (int i = 0; i < output.x; i++) {
3       for (int j = 0; j < output.y; j++) {
4           double accumulator = 0.0;
5           int iOffset = 0;
6           int jOffset = 0;
7           for (int k = 0; k < convFilter.values.length; k++) {
8               final int newIdx = ((i+iOffset) * y) + (j+jOffset);
9               accumulator = Math.fma(values[newIdx], convFilter.values[k], accumulator);
10              jOffset++;
11              if (jOffset == convFilter.y) {
12                  iOffset++;
13                  jOffset = 0;
14              }
15          }
16          output.add(i,j,accumulator);
17      }
18  }
```

# SIMD convolution

```java
1  static final VectorSpecies<Double> D_S = DoubleVector.SPECIES_PREFERRED;
2
3  // Compute convolution
4  for (int i = 0; i < output.x; i++) {
5      for (int j = 0; j < output.y; j++) {
6          fillBuffer(buffer, i, j, convFilter.x, convFilter.y);
7          int k = 0;
8          var accumVec = DoubleVector.zero(D_S);
9          for (; k < D_S.loopBound(buffer.length); k += D_S.length()) {
10             var data = DoubleVector.fromArray(D_S, buffer, k);
11             var conv = DoubleVector.fromArray(D_S, convFilter.values, k);
12             accumVec = data.fma(conv,accumVec);
13         }
14
15         double accumulator = accumVec.reduceLanes(VectorOperators.ADD);
16         for (; k < buffer.length; k++) {
17             accumulator = Math.fma(buffer[k], convFilter.values[k], accumulator);
18         }
19         output.add(i, j, accumulator);
20     }
21 }
```

# Performance is ok (higher is better)

```
// On 8 core VM using an Intel Xeon Platinum 8358 (Ice Lake, AVX-512)
Benchmark                                  Mode   Cnt   Score     Error    Units
Benchmarks.scalarBufferedConv              thrpt    5   0.245 ±   0.002  ops/ms
Benchmarks.scalarConv                      thrpt    5   0.428 ±   0.001  ops/ms
Benchmarks.simdBufferedConv                thrpt    5   0.459 ±   0.002  ops/ms
```

- Performance is not optimal due to the copy needed to line up the image

- The buffered SIMD implementation is twice as fast as the buffered scalar one
  - But only slightly faster than the unbuffered scalar implementation with more complex indexing

- If we know the convolution size we can hand craft code for just that operation

- Nvidia do this in cuDNN (their GPU machine learning primitives library)

- With the Vector API, we can do the same thing in Java

# Specialized 7x7 SIMD convolution

```java
1   // Load filter into registers
2   final boolean[] sevenLaneMask = new boolean[]{true,true,true,true,true,true,true,false};
3   final var mask = VectorMask.fromArray(D_S, sevenLaneMask, 0);
4   final var one = DoubleVector.fromArray(D_S, convFilter.values, 0, mask);
5   final var two = DoubleVector.fromArray(D_S, convFilter.values, 7, mask);
6   ...
7   final var seven = DoubleVector.fromArray(D_S, convFilter.values, 42, mask);
8
9   // Compute convolution
10  for (int i = 0; i < output.x; i++) {
11      for (int j = 0; j < output.y; j++) {
12          var outputVec = DoubleVector.zero(D_S);
13
14          var input = DoubleVector.fromArray(D_S, values, (i*y) + j);
15          outputVec = input.fma(one, outputVec);
16          input = DoubleVector.fromArray(D_S, values, ((i+1)*y) + j);
17          outputVec = input.fma(two, outputVec);
18          ...
19          input = DoubleVector.fromArray(D_S, values, ((i+6)*y) + j, mask);
20          outputVec = input.fma(seven, outputVec);
21
22          output.add(i, j, outputVec.reduceLanes(VectorOperators.ADD));
23      }
24  }
```

# Performance is improved

```
// On 8 core VM using an Intel Xeon Platinum 8358 (Ice Lake, AVX-512)
Benchmark                                   Mode  Cnt   Score     Error   Units
Benchmarks.scalarBufferedConv               thrpt   5   0.245 ±   0.002  ops/ms
Benchmarks.scalarConv                       thrpt   5   0.428 ±   0.001  ops/ms
Benchmarks.simdBufferedConv                 thrpt   5   0.459 ±   0.002  ops/ms
Benchmarks.simdSpecializedConv              thrpt   5   1.379 ±   0.012  ops/ms
```

- The loop is unrolled and specialized to convolution size & vector width

- We lose 1/8th of our computation each time as the 8th lane is masked

  — But C2 can promote the convolution matrix into vector registers

- So we load the convolution matrix once, and reuse it for each pixel

  — Resulting in a 3x faster implementation, at the cost of more complicated code

# Example: lanewise arc-tangent

- The JDK bundles some of the routines from Intel's Short Vector Math Library for lanewise transcendental operations

  — Same accuracy as the scalar methods in `java.lang.Math`

- We can increase performance for reduced accuracy using a different implementation

- A good implementation that balances performance with accuracy can be found in the Cephes Mathematical Library

  — https://github.com/jeremybarnes/cephes/blob/master/single/atanf.c

# Cephes atanf algorithm pseudo code

```
1   // Range reduction to interval [0, tan(pi/8)]
2   neg = false
3   if (x < 0) {
4       neg = true
5       x = -x
6   }
7
8   A0 = 0
9   if (x > tan(3pi/8)) {
10      x = -(1 / x)
11      A0 = pi/2
12  } else (if x > tan(pi/8)) {
13      x = (x - 1) / (x + 1)
14      A0 = pi/4
15  }
16  y = A9 * x^9 - A7 * x^7 + A5 * x^5 - A3 * x^3 + x + A0
17
18  if (neg) y = -y
```

# Absolute value and extract the sign bit

```
1  IntVector sign_bit = x.reinterpretAsInts();
2  // Take the absolute value
3  x = sign_bit.and(INV_SIGN_MASK_I).reinterpretAsFloats();
4  // Extract the sign bit (upper one) to be applied to the result
5  sign_bit = sign_bit.and(SIGN_MASK_I);
```

# Comparisons producing masks for blending

```java
1  // x > tan(3pi/8)
2  VectorMask<Float> cmpHi = x.compare(VectorOperators.GT, ATANF_RANGE_HI_F);
3  // x > tan(pi/8)
4  VectorMask<Float> cmpLow = x.compare(VectorOperators.GT, ATANF_RANGE_LOW_F);
6  // x > tan(pi/8) && !( x > tan(3pi/8) )
7  VectorMask<Float> cmpLowHi = cmpLow.andNot(cmpHi);
```

# Blend in x and A0 according to range of x

```
1   // x0 = -1.0 / x
2   FloatVector x0 = FloatVector.broadcast(F_S, -1.0f).div(x);
3   // x1 = (x - 1.0) / (x + 1.0)
4   FloatVector x1 = x.sub(1.0f).div(x.add(1.0f));
5
6   // Blend in -1 / x and (x - 1) / (x + 1) and x to x
7   // x0 if x > tan(3pi/8)
8   x = x.blend(x0, cmpHi);
9   // x1 if x > tan(pi/8) && !( x > tan(3pi/8) )
11  x = x.blend(x1, cmpLowHi);
12
13  // Blend in pi/2 and pi/4 and 0 to A0 coefficient
14  // pi/2 if x > tan(3pi/8)
15  FloatVector y0 = FloatVector.zero(F_S).blend(HALF_PI_F, cmpHi);
16  // pi/4 if x > tan(pi/8) && !( x > tan(3pi/8) )
17  FloatVector y = y0.blend(QUATER_PI_F, cmpLowHi);
```

# Compute the polynomial

```
1   // A9 * x^9 − A7 * x^7 + A5 * x^5 − A3 * x^3 + x + A0
2   FloatVector zz = x.mul(x);
3   FloatVector acc = zz.fma(
4           FloatVector.broadcast(F_S, ATANF_COF_P0_F),
5           FloatVector.broadcast(F_S, −ATANF_COF_P1_F));
6   acc = acc.fma(
7           zz,
8           FloatVector.broadcast(F_S, ATANF_COF_P2_F));
9   acc = acc.fma(
10          zz,
11          FloatVector.broadcast(F_S, −ATANF_COF_P3_F));
12  acc = acc.mul(zz);
13  acc = acc.fma(x, x);
14  // Add A0 coefficient
15  y = y.add(acc);
```

# Restore the sign

```
1  y = y.reinterpretAsInts()
2          .lanewise(VectorOperators.XOR, sign_bit)
3          .reinterpretAsFloats();
```

# Performance is good (lower is better)

```
// 2.6 GHz 6-Core Intel Core i7
// -XX:-TieredCompilation
Benchmark  (size)  Mode  Cnt      Score       Error  Units
scalar       1024  avgt   10   3437.683 ±   65.388  ns/op
vector       1024  avgt   10    494.722 ±    3.962  ns/op
```

# Guidance

- Measure using a benchmark tool such as JMH
  - Use perfasm (on Linux) or dtraceasm (on Mac) to look at generated code
- Store species in static final fields
  - This will ensure the runtime compiler fold away constant expressions
- Avoid storing vectors in instance fields or arrays
  - Keep use to local variables and method arguments and return values
  - Static fields can be used for vector constants

- Avoid identity sensitive operations on vectors or assigning them to `null`
- Take care when splitting an algorithm into multiple methods
  - Method calls need to inline
  - No support vector calling conventions

# Status

- JEP 426: Vector API (Fourth Incubator) delivered in JDK 19
  - Compile and run with `--add-modules=jdk.incubator.vector`
  - Solid progress made on each delivery to optimize the implementation
- Work has started to align with value classes and types from Project Valhalla
  - Vectors are a special kind of value
- Experimental support for vectors of 16 bit floating point numbers
  - See https://github.com/openjdk/panama-vector/tree/vectorIntrinsics+fp16
  - This also requires value classes and types from Project Valhalla