# Cross-Element Vectorization in Firedrake

TJ Sun (ts2914@ic.ac.uk), Lawrence Mitchell, David A Ham, Paul H J Kelly
June 2018

# loo.py as our new backend

Introducing loo.py
- Andreas Klöckner et al. (UIUC)
- Based on polyhedral model of loops
- ≈ model of loops + transformations
- "Do what I tell you"
- Support multiple backends
  - CPU
  - ISPC
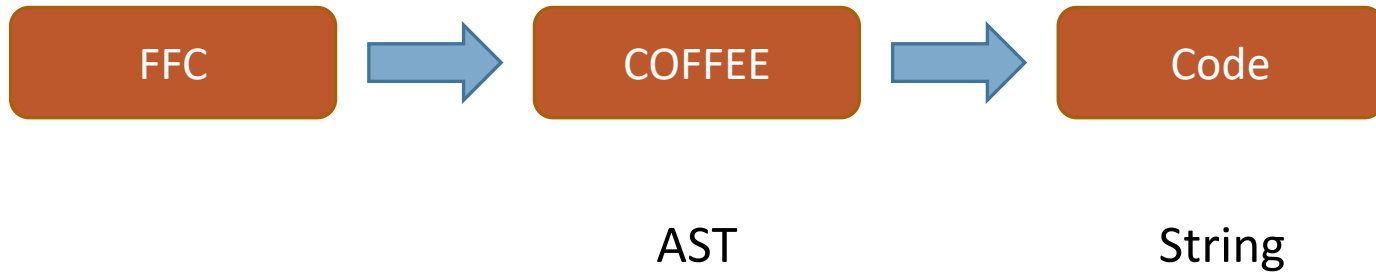  - OpenCL, PyOpenCL
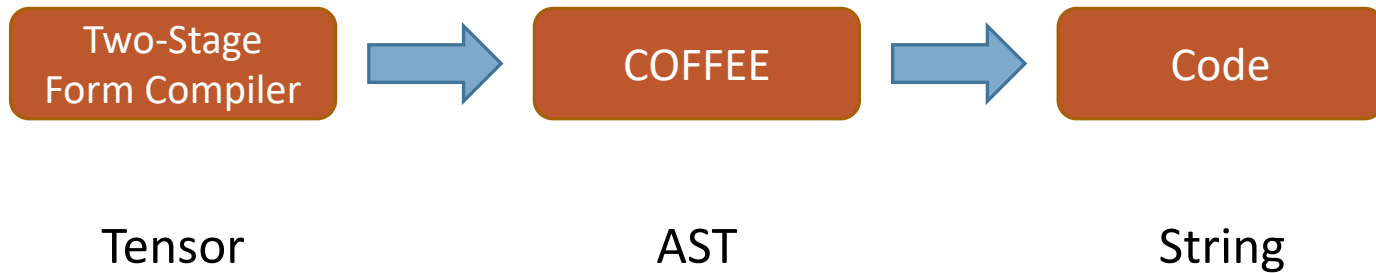  - Cuda
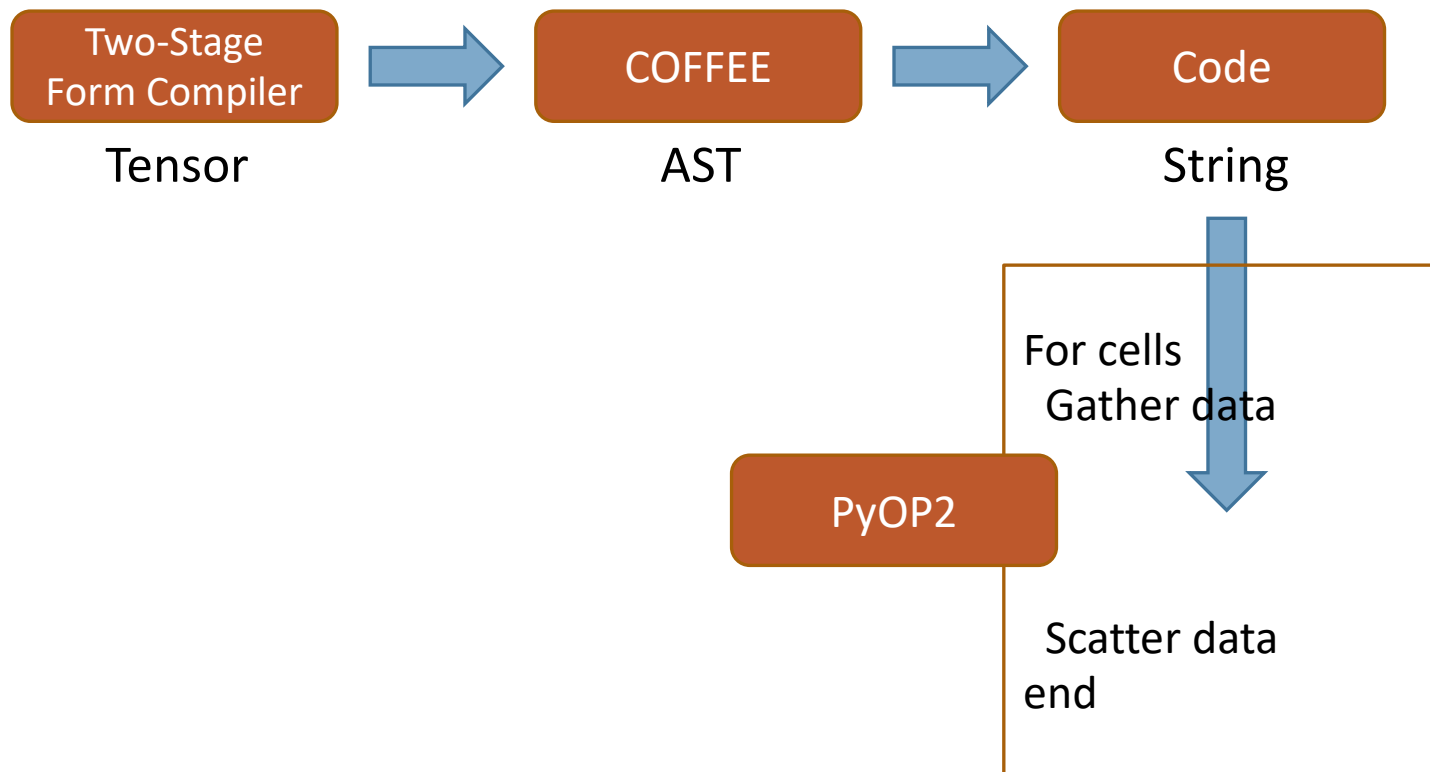
# Code generation for assembly

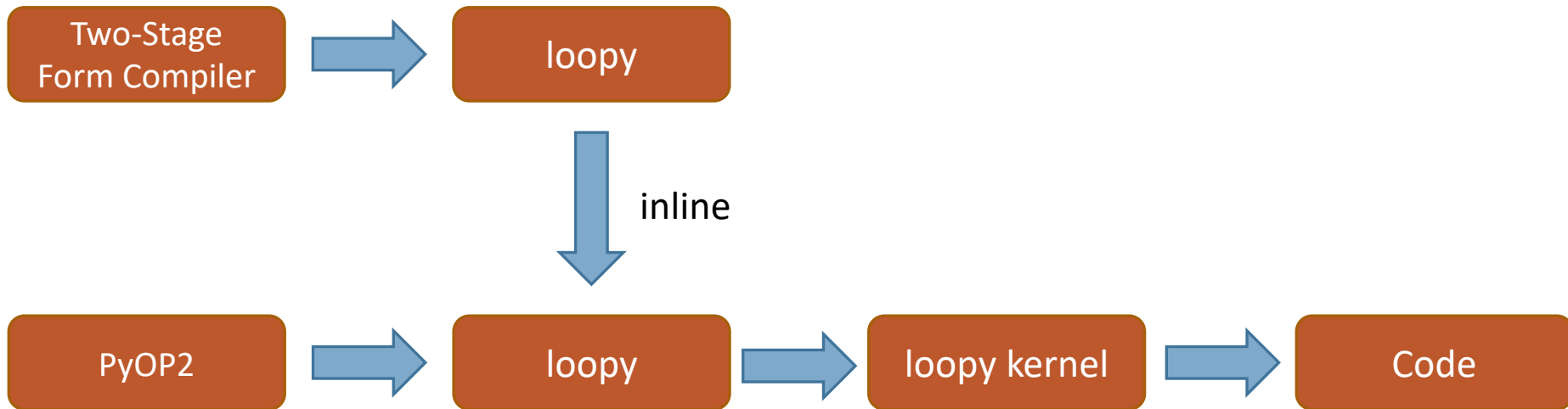FFC → Code

string

# Code generation for assembly

FFC → COFFEE → Code

AST          String

# Code generation for assembly

| Two-Stage Form Compiler | → | COFFEE | → | Code |
|---|---|---|---|---|
| Tensor | | AST | | String |

# Code generation for assembly

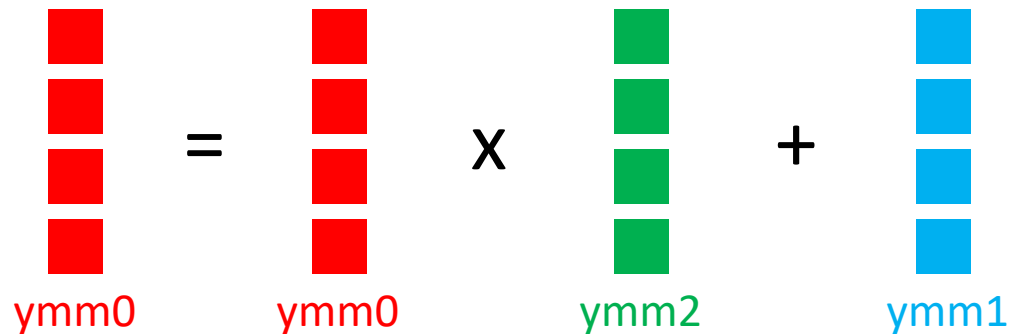| Two-Stage Form Compiler | → | COFFEE | → | Code |
|:---:|:---:|:---:|:---:|:---:|
| Tensor | | AST | | String |

PyOP2

For cells
Gather data

Scatter data
end

# Code generation for assembly

# What is vectorization

- SIMD (single instruction multiple data) programming model

- e.g. VFMADD213PD ymm0, ymm1, ymm2 (in AVX2 instruction set)

ymm0 = ymm0 x ymm2 + ymm1

- 8 double precision operations with 1 instruction

- Need to issue 2 FMA instructions in 1 cycle to get advertised performance

- SIMD width doubles every 4 years
  - AVX512 instructions can do 8 doubles

- Naïve code usually achieves <10% peak performance

```
 3    static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const
 4    {
 5
 6      /* ... */
 7
 8      for (int ip = 0; ip <= 5; ++ip)
 9      {
10
11        /* ... */
12
13        for (int j0 = 0; j0 <= 9; ++j0)
14        {
15          A[2 * j0] = A[2 * j0] + t18[10 * ip + j0] * t36 + t19[10 * ip + j0] * t35;
16          A[1 + 2 * j0] = A[1 + 2 * j0] + t18[10 * ip + j0] * t32 + t19[10 * ip + j0] * t31;
17        }
18      }
19    }
20
21    void wrapper(int const start, int const end, double *__restrict__ dat0, double const *__restrict
22    {
23      /* ... */
24      for (int n = start; n <= -1 + end; ++n)
25      {
26        for (int i7 = 0; i7 <= 2; ++i7)
27          for (int i8 = 0; i8 <= 1; ++i8)
28            t3[2 * i7 + i8] = dat1[2 * map1[3 * n + i7] + i8];
29        for (int i13 = 0; i13 <= 9; ++i13)
30          for (int i14 = 0; i14 <= 1; ++i14)
31            t4[2 * i13 + i14] = dat2[2 * map0[10 * n + i13] + i14];
32        for (int i1 = 0; i1 <= 9; ++i1)
33          for (int i2 = 0; i2 <= 1; ++i2)
34            t2[2 * i1 + i2] = 0.0;
35
36        form0_cell_integral_otherwise(t2, t3, t4);
37
38        for (int i20 = 0; i20 <= 1; ++i20)
39          for (int i19 = 0; i19 <= 9; ++i19)
40            dat0[2 * map0[10 * n + i19] + i20] = dat0[2 * map0[10 * n + i19] + i20] + t2[2 * i19 + i20
41      }
42    }
43
```

Action of linear elasticity operator on triangle mesh, CG3 space

```
  3   static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const
  4   {
  5
  6      /* ... */
  7
  8      for (int ip = 0; ip <= 5; ++ip)
  9      {
 10
 11         /* ... */
 12
 13         for (int j0 = 0; j0 <= 9; ++j0)
 14         {
 15            A[2 * j0] = A[2 * j0] + t18[10 * ip + j0] * t36 + t19[10 * ip + j0] * t35;
 16            A[1 + 2 * j0] = A[1 + 2 * j0] + t18[10 * ip + j0] * t32 + t19[10 * ip + j0] * t31;
 17         }
 18      }
 19   }
 20
 21   void wrapper(int const start, int const end, double *__restrict__ dat0, double const *__restrict
 22   {
 23      /* ... */
 24      for (int n = start; n <= -1 + end; ++n)
 25      {
 26         for (int i7 = 0; i7 <= 2; ++i7)
 27            for (int i8 = 0; i8 <= 1; ++i8)
 28               t3[2 * i7 + i8] = dat1[2 * map1[3 * n + i7] + i8];
 29         for (int i13 = 0; i13 <= 9; ++i13)
 30            for (int i14 = 0; i14 <= 1; ++i14)
 31               t4[2 * i13 + i14] = dat2[2 * map0[10 * n + i13] + i14];
 32         for (int i1 = 0; i1 <= 9; ++i1)
 33            for (int i2 = 0; i2 <= 1; ++i2)
 34               t2[2 * i1 + i2] = 0.0;
 35
 36         form0_cell_integral_otherwise(t2, t3, t4);
 37
 38         for (int i20 = 0; i20 <= 1; ++i20)
 39            for (int i19 = 0; i19 <= 9; ++i19)
 40               dat0[2 * map0[10 * n + i19] + i20] = dat0[2 * map0[10 * n + i19] + i20] + t2[2 * i19 + i20
 41      }
 42   }
 43
```
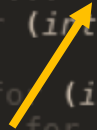
Kernel

Wrapper

Action of linear elasticity operator on triangle mesh, CG3 space

```
3    static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const
4    {
5
6      /* ... */
7
8      for (int ip = 0; ip <= 5; ++ip)
9      {
10
11       /* ... */
12
13       for (int j0 = 0; j0 <= 9; ++j0)
14       {
15         A[2 * j0] =                              ;
16         A[1 + 2 * j0                              0] * t31;
17       }
18     }
19    }
20
21    void wrapper(int const start, int const end, double *__restrict__ dat0, double const *__restrict_
22    {
23      /* ... */
24      for (int n = start; n <= -1 + end; ++n)
25      {
26        for (int i7 = 0; i7 <= 2; ++i7)
27          for (int i8 = 0; i8 <= 1; ++i8)
28            t3[2 * i7 + i8] = dat1[2 * map1[3 * n + i7] + i8];
29        for (int i13 = 0; i13 <= 9; ++i13)
30          for (int i14 = 0; i14 <= 1; ++i14)
31            t4[2 * i13 + i14] = dat2[2 * map0[10 * n + i13] + i14];
32        for (int i1 = 0; i1 <= 9; ++i1)
33          for (int i2 = 0; i2 <= 1; ++i2)
34            t2[2 * i1 + i2] = 0.0;
35
36        form0_cell_integral_otherwise(t2, t3, t4);
37
38        for (int i20 = 0; i20 <= 1; ++i20)
39          for (int i19 = 0; i19 <= 9; ++i19)
40            dat0[2 * map0[10 * n + i19] + i20] = dat0[2 * map0[10 * n + i19] + i20] + t2[2 * i19 + i20
41      }
42    }
43
```

Outer loop over all elements in the mesh

```
for (int n = start; n <= -1 + end; ++n)
```

Wrapper

Action of linear elasticity operator on triangle mesh, CG3 space

```
 3   static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const
 4   {
 5
 6       /* ... */
 7
 8       for (int ip = 0
 9       {
10
11           /*
12
13           for
14           {
15               A
16               A
17           }
18       }
19   }
20
21   void wr
22   {
23       /* ...
24       for (int n = start; n <= -1 + end; ++n)
25       {
26           for (int i7 = 0; i7 <= 2; ++i7)
27               for (int i8 = 0; i8 <= 1; ++i8)
28                   t3[2 * i7 + i8] = dat1[2 * map1[3 * n + i7] + i8];
29           for (int i13 = 0; i13 <= 9; ++i13)
30               for (int i14 = 0; i14 <= 1; ++i14)
31                   t4[2 * i13 + i14] = dat2[2 * map0[10 * n + i13] + i14];
32           for (int i1 = 0; i1 <= 9; ++i1)
33               for (int i2 = 0; i2 <= 1; ++i2)
34                   t2[2 * i1 + i2] = 0.0;
35
36           form0_cell_integral_otherwise(t2, t3, t4);
37
38           for (int i20 = 0; i20 <= 1; ++i20)
39               for (int i19 = 0; i19 <= 9; ++i19)
40                   dat0[2 * map0[10 * n + i19] + i20] = dat0[2 * map0[10 * n + i19] + i20] + t2[2 * i19 + i20
41       }
42   }
43
```

**Indirect gathering of input data for kernel**

```
for (int i7 = 0; i7 <= 2; ++i7)
    for (int i8 = 0; i8 <= 1; ++i8)
        t3[2 * i7 + i8] = dat1[2 * map1[3 * n + i7] + i8];
for (int i13 = 0; i13 <= 9; ++i13)
    for (int i14 = 0; i14 <= 1; ++i14)
        t4[2 * i13 + i14] = dat2[2 * map0[10 * n + i13] + i14];
for (int i1 = 0; i1 <= 9; ++i1)
    for (int i2 = 0; i2 <= 1; ++i2)
        t2[2 * i1 + i2] = 0.0;
```

Action of linear elasticity operator on triangle mesh, CG3 space

```
 3  static inline void kernel(double * __restrict__ A, double const * __restrict__ coords, double const
 4  {
 5
 6    /* ... */
 7
 8    for (int ip = 0; ip <= 5; ++ip)
 9    {
10
11      /* ... */
12
13      for (int j0 = 0; j0 <= 9; ++j0)
14      {
15        A[2 * j0] = A[2 * j0] + t18[10 * ip + j0] * t36 + t19[10 * ip + j0] * t35;
16        A[1 + 2 * j0] = A[1 + 2 * j0] + t18[10 * ip + j0] * t32 + t19[10 * ip + j0] * t31;
17      }
18    }
19  }
20
21  void wrapper(int const start, int const end, double * __restrict__ dat0, double const * __restrict__
22  {
23    /* ... */
24    for (int n = start; n <= -1 + en
25    {
26      for (int i7 = 0; i7 <= 2; ++i7
27        for (int i8 = 0; i8 <= 1; ++i8)
28          t3[2 * i7 + i8] = dat1[2 *
29      for (int i13 = 0; i13 <= 9; ++
30        for (int i14 = 0; i14 <= 1;
31          t4[2 * i13 + i14] = dat2[2
32      for (int i1 = 0; i1 <= 9; ++i1)
33        for (int i2 = 0; i2 <= 1; ++i2)
34          t2[2 * i1 + i2] = 0.0;
35
36      form0_cell_integral_otherwise(t2, t3, t4);
37
38      for (int i20 = 0; i20 <= 1; ++i20)
39        for (int i19 = 0; i19 <= 9; ++i19)
40          dat0[2 * map0[10 * n + i19] + i20] = dat0[2 * map0[10 * n + i19] + i20] + t2[2 * i19 + i20
41    }
42  }
43
```

Kernel "call", actually it is inlined

```
kernel(t2, t3, t4);
```

Action of linear elasticity operator on triangle mesh, CG3 space

```
 3  static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const
 4  {
 5
 6    /* ... */
 7
 8    for (int ip = 0; ip <= 5; ++ip)
 9    {
10
11      /* ... */
12
13      for (int j0 = 0; j0 <= 9; ++j0)
14      {
15        A[2 * j0] = A[2 * j0] + t18[10 * ip + j0] * t36 + t19[10 * ip + j0] * t35;
16        A[1 + 2 * j0] = A[1 + 2 * j0] + t18[10 * ip + j0] * t32 + t19[10 * ip + j0] * t31;
17      }
18    }
19  }
20
21  void wrapper(int const start, int const end, double *__restrict__ dat0, double const *__restrict_
22  {
23    /* ... */
24    for (int n = start; n <= -1 + end; ++n)
25    {
26      for (int i
27        for (int
28          t3[2
29
30    for (int i20 = 0; i20 <= 1; ++i20)
31      for (int i19 = 0; i19 <= 9; ++i19)
32        dat0[2 * map0[10 * n + i19] + i20] += t2[2 * i19 + i20];
33
34
35
36    form0_cell_integral_otherwise(t2, t3, t4);
37
38    for (int i20 = 0; i20 <= 1; ++i20)
39      for (int i19 = 0; i19 <= 9; ++i19)
40        dat0[2 * map0[10 * n + i19] + i20] = dat0[2 * map0[10 * n + i19] + i20] + t2[2 * i19 + i20
41  }
42  }
43
```

Indirect scattering of local tensor to global tensor

Action of linear elasticity operator on triangle mesh, CG3 space

```
  3     static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const
  4     {
  5
  6       /* ... */
  7
  8       for (int ip = 0; ip <= 5; ++ip)
  9       {
 10
 11         /* ... */
 12
 13         for (int j0 = 0; j0 <= 9; ++j0)
 14         {
 15           A[2 * j0] = A[2 * j0] + t18[10 * ip + j0] * t36 + t19[10 * ip + j0] * t35;
 16           A[1 + 2 * j0] = A[1 + 2 * j0] + t18[10 * ip + j0] * t32 + t19[10 * ip + j0] * t31;
 17         }
 18       }
 19     }
 20
 21     void wrapper
 22     for (int ip = 0; ip <= 5; ++ip)
 23     {
 24
 25
 26       /*  ...
 27
 28
 29       for (int j0 = 0; j0 <= 9; ++j0)
 30       {
 31
 32         A[2 * j0] += t18[10 * ip + j0] * t36 + t19[10 * ip + j0] * t35;
 33         A[1 + 2 * j0] += t18[10 * ip + j0] * t32 + t19[10 * ip + j0] * t31;
 34       }
 35     }
 36
 37
 38       for (int i20 = 0; i20 <= 1; ++i20)
 39         for (int i19 = 0; i19 <= 9; ++i19)
 40           dat0[2 * map0[10 * n + i19] + i20] = dat0[2 * map0[10 * n + i19] + i20] + t2[2 * i19 + i20
 41       }
 42     }
 43
```

Kernel

Outer loop: contraction over quadrature points

inner loop over degrees of freedom

Action of linear elasticity operator on triangle mesh, CG3 space

# Vectorization strategy

o "Intra-kernel" can be tricky
- Trip count can be small and/or not multiple of SIMD width
- Alignment to cache boundary
- Stride 1 access
- Operations outside of innermost loop not vectorized
- Loop structure varies with PDE, discretization, mesh
- And we have done many of these in firedrake [1]

o "Inter-kernel" provides a generic solution
- Vector-expand the kernel to act on N elements together, N=SIMD width
- All operations can be vectorized
- Can always do this systematically
- Downside: increasing working size

[1] F. Luporini, et al.  ACM TACO 2015

```
 3    static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const *__restrict
 4    {
 5
 6        /* ... */
 7
 8        for (int ip = 0; ip <= 5; ++ip)
 9        {
10            /* ... */
11
12            for (int j0 = 0; j0 <= 9; ++j0)
13                for (int elem = 0; elem <= 3; ++elem)
14                {
15                    A[elem + 4 * 2 * j0] = A[elem + 4 * 2 * j0] + t18[10 * ip + j0] * t36[elem] + t19[10 * ip + j0] * t35[
16                    A[elem + 4 * (1 + 2 * j0)] = A[elem + 4 * (1 + 2 * j0)] + t18[10 * ip + j0] * t32[elem] + t19[10 * ip
17                }
18        }
19    }
20
21    void wrap_form0_cell_integral_otherwise(int const start, int const end, double *__restrict__ dat0, double cons
22    {
23
24        /* ... */
25
26        for (int n_outer = (start / 4); n_outer <= ((-4 + end) / 4); ++n_outer)
27        {
28            for (int i7 = 0; i7 <= 2; ++i7)
29                for (int i8 = 0; i8 <= 1; ++i8)
30                    for (int n_inner = 0; n_inner <= 3; ++n_inner)
31                        t3[8 * i7 + 4 * i8 + n_inner] = dat1[2 * map1[3 * (4 * n_outer + n_inner) + i7] + i8];
32            for (int i13 = 0; i13 <= 9; ++i13)
33                for (int i14 = 0; i14 <= 1; ++i14)
34                    for (int n_inner = 0; n_inner <= 3; ++n_inner)
35                        t4[8 * i13 + 4 * i14 + n_inner] = dat2[2 * map0[10 * (4 * n_outer + n_inner) + i13] + i14];
36            for (int i1 = 0; i1 <= 9; ++i1)
37                for (int i2 = 0; i2 <= 1; ++i2)
38                    for (int n_inner = 0; n_inner <= 3; ++n_inner)
39                        t2[8 * i1 + 4 * i2 + n_inner] = 0.0;
40
41            kernel(t2, t3, t4);
42
43            for (int i20 = 0; i20 <= 1; ++i20)
44                for (int i19 = 0; i19 <= 9; ++i19)
45                    for (int n_inner = 0; n_inner <= 3; ++n_inner)
46                        dat0[2 * map0[10 * (4 * n_outer + n_inner) + i19] + i20] = dat0[2 * map0[10 * (4 * n_outer + n_inner
47        }
48    }
```

Action of linear elasticity operator on triangle mesh, batched by 4

```
  3   static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const *__restrict
  4   {
  5
  6     /* ... */
  7
  8     for (int ip = 0; ip <= 5; ++ip)
  9     {
 10       /* ... */
 11
 12       for (int j0 =
 13         for (int ele
 14         {
 15           A[elem + 4                                                    36[elem] + t19[10 * ip + j0] * t35[
 16           A[elem + 4                                                    ip + j0] * t32[elem] + t19[10 * ip
 17         }
 18   for (int n_outer = (start / 4); n_outer <= ((-4 + end) / 4); ++n_outer)
 19   }
 20
 21   void wrap_form0_cell_integral_otherwise(int const start, int const end, double *__restrict__ dat0, double cons
 22   {
 23
 24     /* ... */
 25
 26     for (int n_outer = (start / 4); n_outer <= ((-4 + end) / 4); ++n_outer)
 27     {
 28       for (int i7 = 0; i7 <= 2; ++i7)
 29         for (int i8 = 0; i8 <= 1; ++i8)
 30           for (int n_inner = 0; n_inner <= 3; ++n_inner)
 31             t3[8 * i7 + 4 * i8 + n_inner] = dat1[2 * map1[3 * (4 * n_outer + n_inner)
 32       for (int i13 = 0; i13 <= 9; ++i13)
 33         for (int i14 = 0; i14 <= 1; ++i14)
 34           for (int n_inner = 0; n_inner <= 3; ++n_inner)
 35             t4[8 * i13 + 4 * i14 + n_inner] = dat2[2 * map0[10 * (4 * n_outer + n_inner) + i13] + i14];
 36       for (int i1 = 0; i1 <= 9; ++i1)
 37         for (int i2 = 0; i2 <= 1; ++i2)
 38           for (int n_inner = 0; n_inner <= 3; ++n_inner)
 39             t2[8 * i1 + 4 * i2 + n_inner] = 0.0;
 40
 41       kernel(t2, t3, t4);
 42
 43       for (int i20 = 0; i20 <= 1; ++i20)
 44         for (int i19 = 0; i19 <= 9; ++i19)
 45           for (int n_inner = 0; n_inner <= 3; ++n_inner)
 46             dat0[2 * map0[10 * (4 * n_outer + n_inner) + i19] + i20] = dat0[2 * map0[10 * (4 * n_outer + n_inner
 47   }
 48   }
```

Split n into n_outer and n_inner
Outer loop stride 4

Wrapper

Action of linear elasticity operator on triangle mesh, batched by 4

Action of linear elasticity operator on triangle mesh, batched by 4

```
3    static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const *__restrict_
4    {
5
6      /* ... */
7
8      for (int ip = 0; ip <= 5; ++ip)
9      {
10       /* ... */
11
12       for (int j0 = 0; j0 <= 9; ++j0)
13         for (int elem = 0; elem <= 3; ++elem)
14         {
15           A[elem + 4 * 2 * j0] = A[elem + 4 * 2 * j0] + t18[10 * ip + j0] * t36[elem] + t19[10 * ip + j0] * t35[
16           A[elem + 4 * (1 + 2 * j0)] = A[elem + 4 * (1 + 2 * j0)] + t18[10 * ip + j0] * t32[elem] + t19[10 * ip
17         }
18     }
19   }
20
21   void wrap_form0_cell_integral_otherwise(int const start, int const end, double *__restrict__ dat0, double cons
22   {
23
24     /* ... */
25
26     for (int n_outer = (start / 4); n_outer <= ((-4 + end) / 4); ++n_outer)
27     {
28       for (int i7 = 0; i7 <= 2; ++i7)
29         for (int i8 = 0; i8 <= 1;
30           for (int n_inner = 0; n
31             t3[8 * i7 + 4 * i8 +                                      outer + n_inner) + i7] + i8];
32       for (int i13 = 0; i13
33         for (int i14 = 0;
34           for (int n_inner
35             t4[8 * i13 +                                             inner) + i13] + i14];
36       for (int i1 = 0; i1
37         for (int i2 =
38           for (int n_inner = 0; n_inner <= 3; ++n_inner)
39             t2[8 * i1 + 4 * i2 + n_inner] = 0.0;
40
41       kernel(t2, t3, t4);
42
43       for (int i20 = 0; i20 <= 1; ++i20)
44         for (int i19 = 0; i19 <= 9; ++i19)
45           for (int n_inner = 0; n_inner <= 3; ++n_inner)
46             dat0[2 * map0[10 * (4 * n_outer + n_inner) + i19] + i20] = dat0[2 * map0[10 * (4 * n_outer + n_inner
47     }
48   }
```

Kernel call

```
kernel(t2, t3, t4);
```

Action of linear elasticity operator on triangle mesh, batched by 4

```
3    static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const *__restrict__
4    {
5
6      /* ... */
7
8      for (int ip = 0; ip <= 5; ++ip)
9      {
10       /* ... */
11
12       for (int j0 = 0; j0 <= 9; ++j0)
13         for (int elem = 0; elem <= 3; ++elem)
14         {
15           A[elem + 4 * 2 * j0] = A[elem + 4 * 2 * j0] + t18[10 * ip + j0] * t36[elem] + t19[10 * ip + j0] * t35[
16           A[elem + 4 * (1 + 2 * j0)] = A[elem + 4 * (1 + 2 * j0)] + t18[10 * ip + j0] * t32[elem] + t19[10 * ip
17         }
18     }
19   }
20
21   void wrap_form0_cell_integral_otherwise(int const start, int const end, double *__restrict__ dat0, double cons
22   {
23
24     /* ... */
25
26     for (int n_outer = (start / 4); n_outer <= ((-4 + end) / 4); ++n_outer)
27     {
28       for (int i7 = 0; i7 <= 2; ++i7)
29         for (int i8 = 0; i8
30           for (int n_inner
31             t3[8 * i7 + 4 *                                                           7] + i8];
```

## Scattering might have race condition

```
for (int i20 = 0; i20 <= 1; ++i20)
  for (int i19 = 0; i19 <= 9; ++i19)
    for (int n_inner = 0; n_inner <= 3; ++n_inner)
      dat0[2 * map0[10 * (4 * n_outer + n_inner) + i19] + i20] += t2[8 * i19 + 4 * i20 + n_inner];
```

```
37         for (int i2 = 0; i2 <= 1; ++i2)
38           for (int n_inner = 0; n_inner <= 3; ++n_inner)
39             t2[8 * i1 + 4 * i2 + n_inner] = 0.0;
40
41     kernel(t2, t3, t4);
42
43     for (int i20 = 0; i20 <= 1; ++i20)
44       for (int i19 = 0; i19 <= 9; ++i19)
45         for (int n_inner = 0; n_inner <= 3; ++n_inner)
46           dat0[2 * map0[10 * (4 * n_outer + n_inner) + i19] + i20] = dat0[2 * map0[10 * (4 * n_outer + n_inner
47   }
48 }
```

Action of linear elasticity operator on triangle mesh, batched by 4

```
static inline void kernel(double *__restrict__ A, double const *__restrict__ coords, double const *__restric
{

  /* ... */

  for (int ip = 0; ip <= 5; ++ip)
  {
    /* ... */

    for (int j0 = 0; j0 <= 9; ++j0)
      for (int elem = 0; elem <= 3; ++elem)
      {
        A[elem + 4 * 2 * j0] = A[elem + 4 * 2 * j0] + t18[10 * ip + j0] * t36[elem] + t19[10 * ip + j0] * t35[
        A[elem + 4 * (1 + 2 * j0)] = A[elem + 4 * (1 + 2 * j0)] + t18[10 * ip + j0] * t32[elem] + t19[10 * ip
      }
  }
}
```

**Kernel**

```
for (int ip = 0; ip <= 5; ++ip)
{
  /* ... */

  for (int j0 = 0; j0 <= 9; ++j0)
    for (int elem = 0; elem <= 3; ++elem)
    {
      A[elem + 4 * 2 * j0] += t18[10 * ip + j0] * t36[elem] + t19[10 * ip + j0] * t35[elem];
      A[elem + 4 * (1 + 2 * j0)] += t18[10 * ip + j0] * t32[elem] + t19[10 * ip + j0] * t31[elem];
    }
}
```

"element" loop pushed to innermost

Trip count 4, stride 1, aligned, independent

```
          for (int n_inner =
            t4[8 * i13 + 4 * i14 + n_inner] = dat2[2 * map0[10 * (4 * n_outer + n_inner) + i13] + i14];
        for (int i1 = 0; i1 <= 9; ++i1)
          for (int i2 = 0; i2 <= 1; ++i2)
            for (int n_inner = 0; n_inner <= 3; ++n_inner)
              t2[8 * i1 + 4 * i2 + n_inner] = 0.0;

      kernel(t2, t3, t4);

        for (int i20 = 0; i20 <= 1; ++i20)
          for (int i19 = 0; i19 <= 9; ++i19)
            for (int n_inner = 0; n_inner <= 3; ++n_inner)
              dat0[2 * map0[10 * (4 * n_outer + n_inner) + i19] + i20] = dat0[2 * map0[10 * (4 * n_outer + n_inner
      }
  }
```
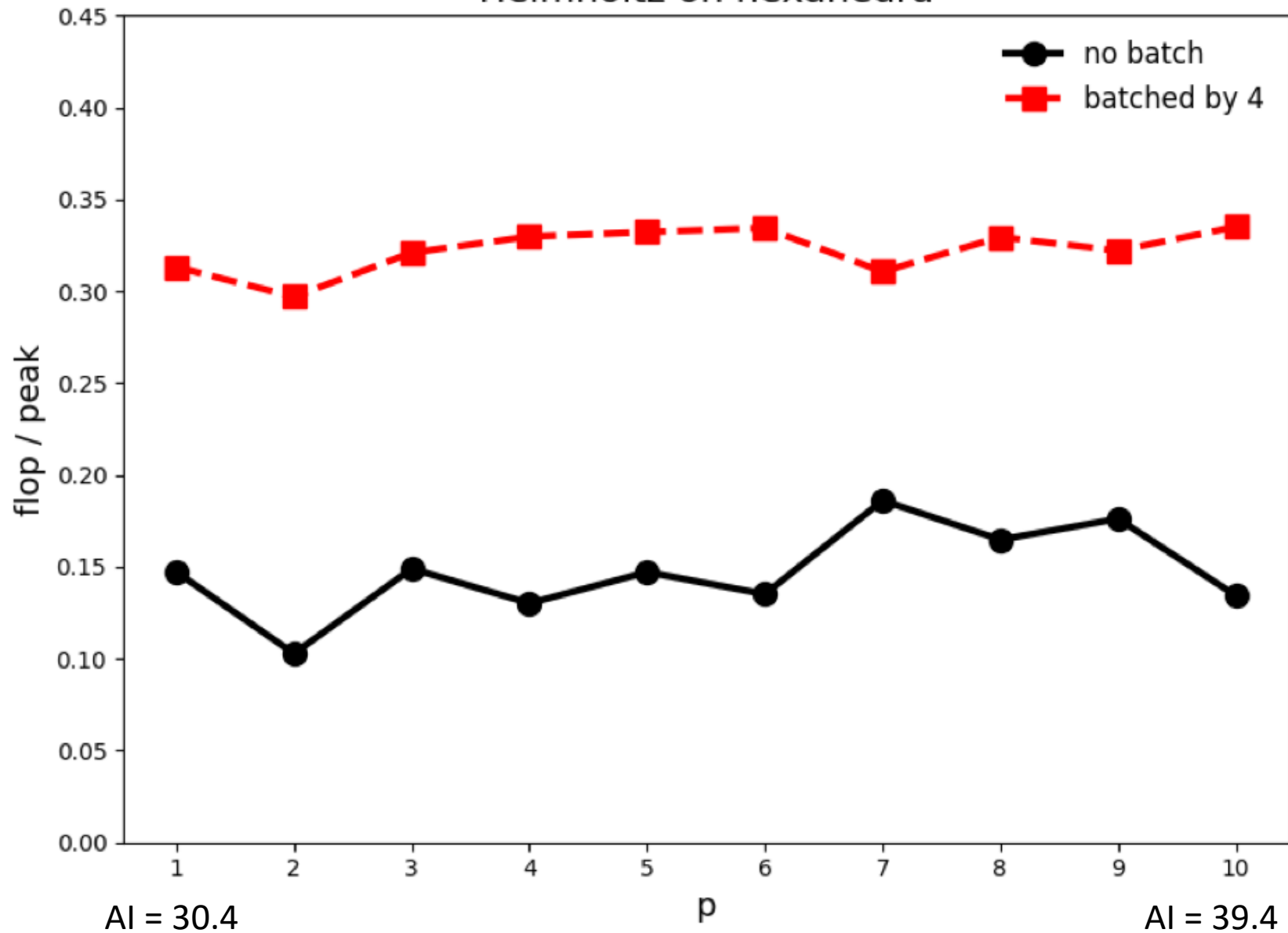
Action of linear elasticity operator on triangle mesh, batched by 4
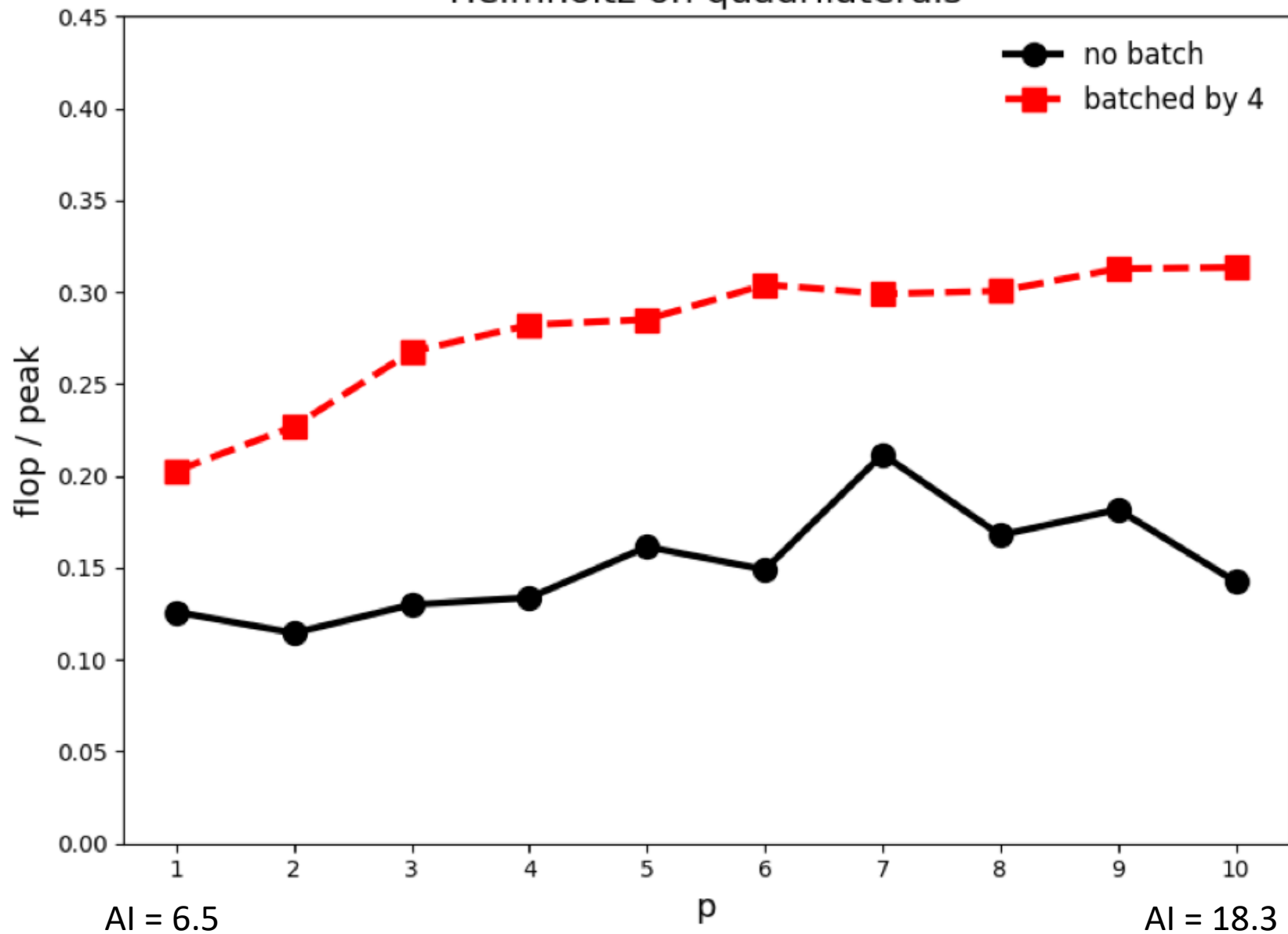
# Experimental setup

o Hardware: Haswell i7-4790 (single core measurement)
  - o SIMD width = 4 (avx2)
  - o Peak flop = 3.6 GHz x 4 (avx2) x 2 (fma) x 2 (issue) = 57.6 Gflops
  - o Running Intel LINPACK binary: 51.0 Gflops
  - o STREAM triad bandwidth: 10.4 GB / s
  - o Roofline AI "regime switching point" = 5.54 flops / byte

o Mesh: hexahedra (3D tensor product element)
  - o TSFC automates sum factorisation [2]
    -> Innermost loop trip count = polynomial degree + 1

o Action of Helmholtz operator
  - o Arithmetic Intensity (perfect cache) 30.4 to 33.6 $\rightarrow$ compute bound

o We present achieved flops / 57.6 Gflops
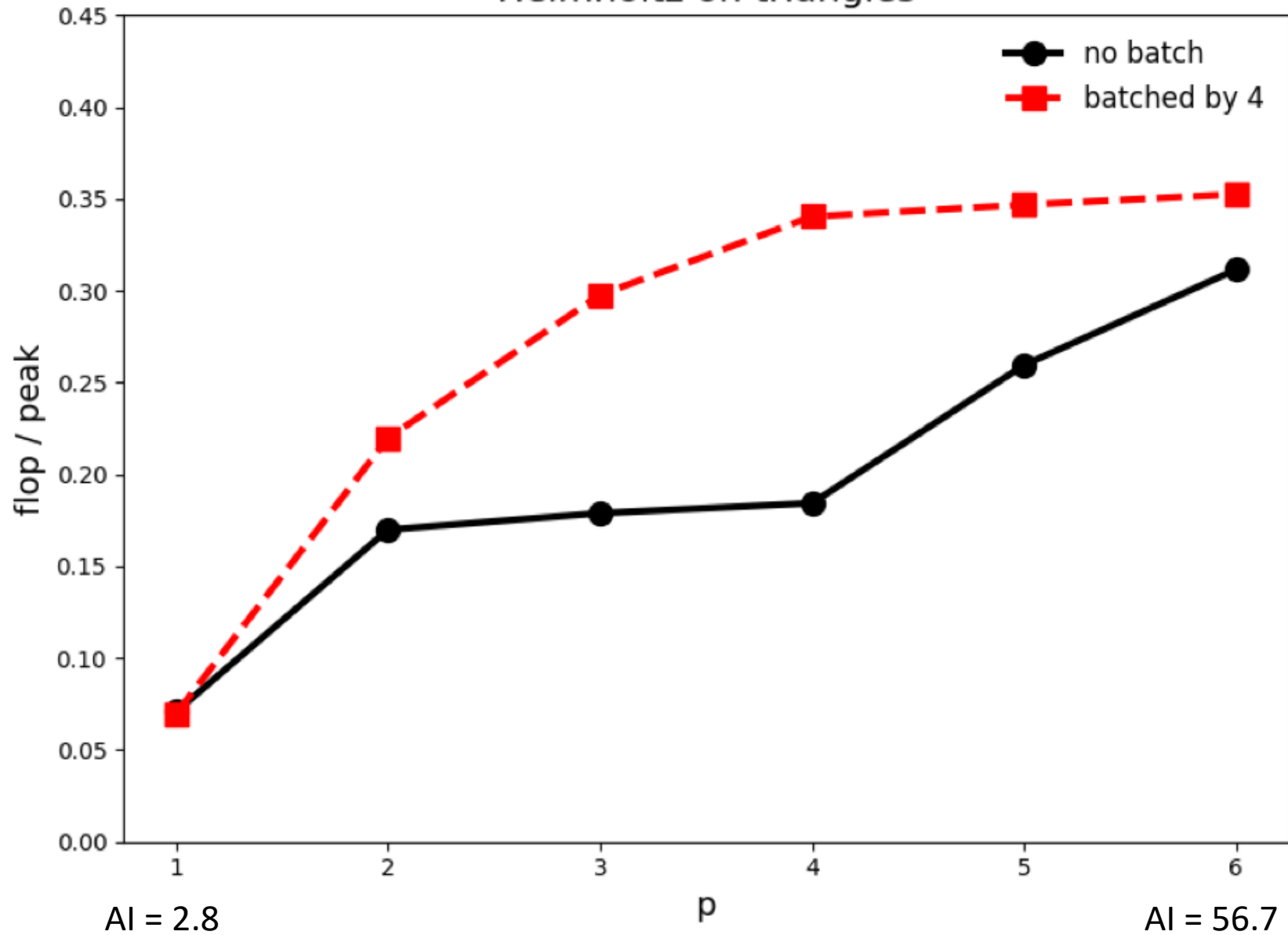
[2] M. Homolya, et al. arXiv:1711.02473 (2017).

Helmholtz on hexahedra

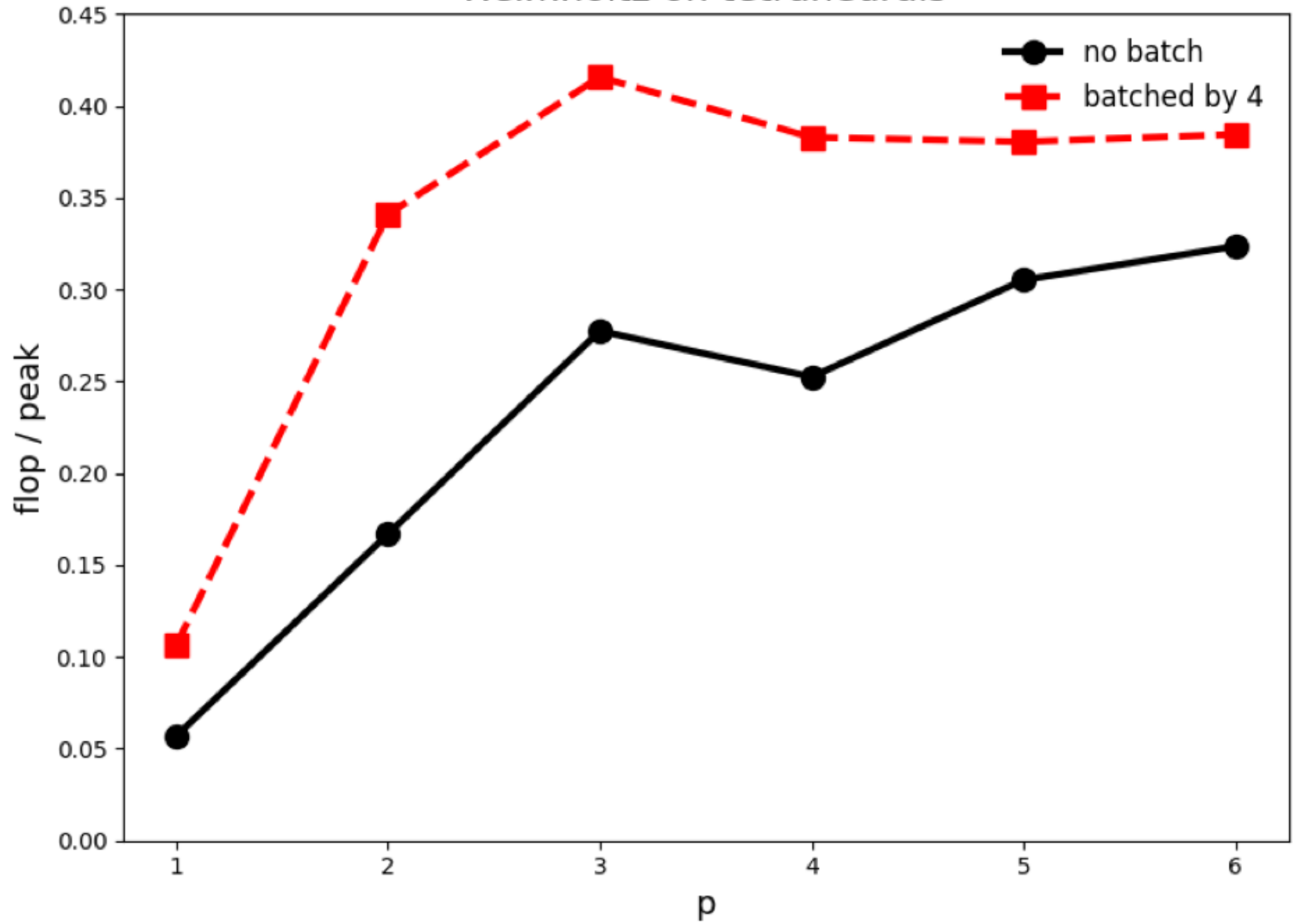AI = 30.4                                                                 AI = 39.4
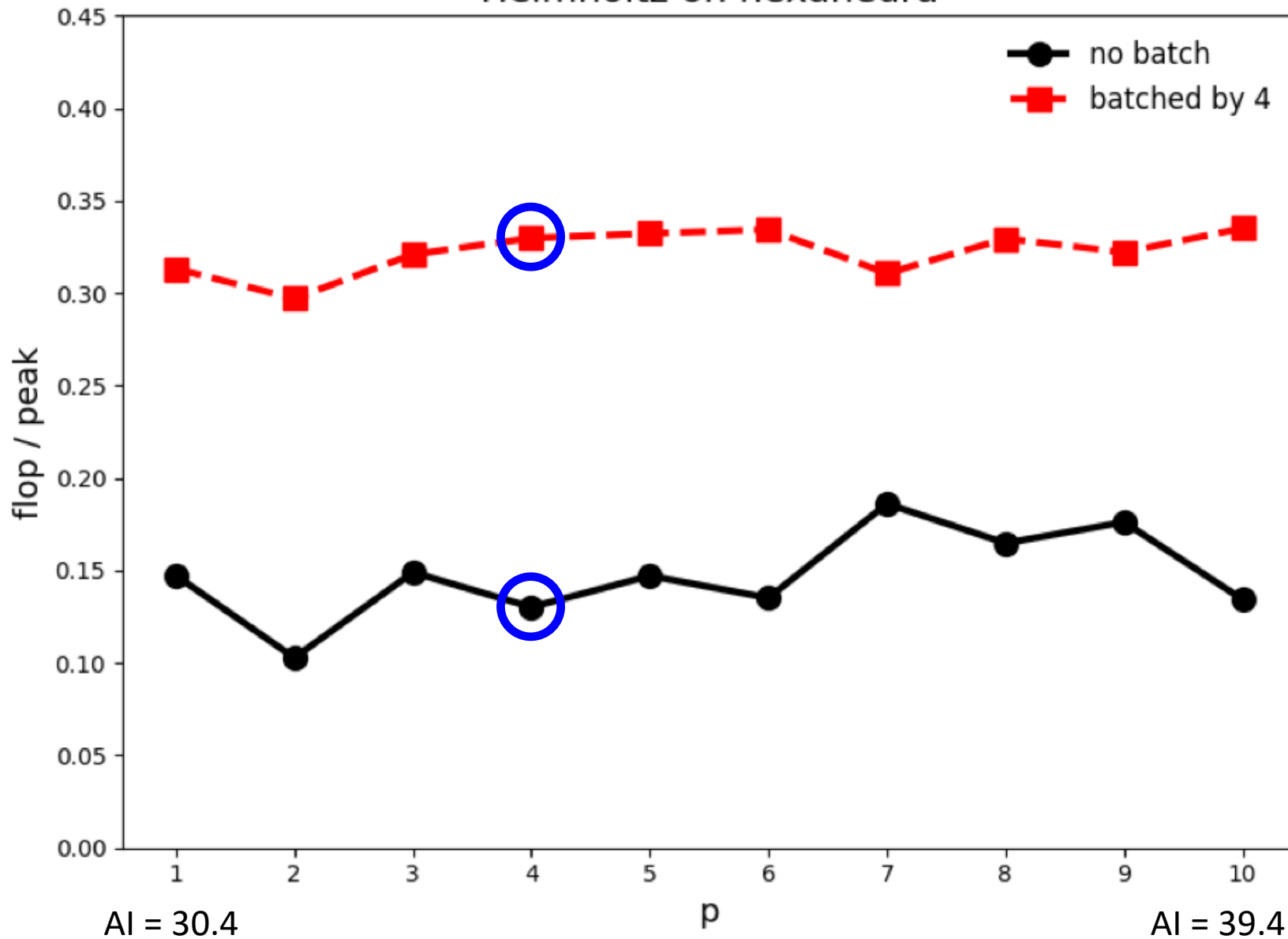
**Helmholtz on quadrilaterals**

AI = 6.5

AI = 18.3

Helmholtz on triangles

AI = 2.8                                          AI = 56.7

Helmholtz on tetrahedrals

Helmholtz on hexahedra

AI = 30.4                    AI = 39.4

# Flop contributions by instruction types

■ AVX2 (4 doubles)   ■ SSE (2 doubles)   ■ scalar (1 double)

Only 0.4% of flops not vectorized

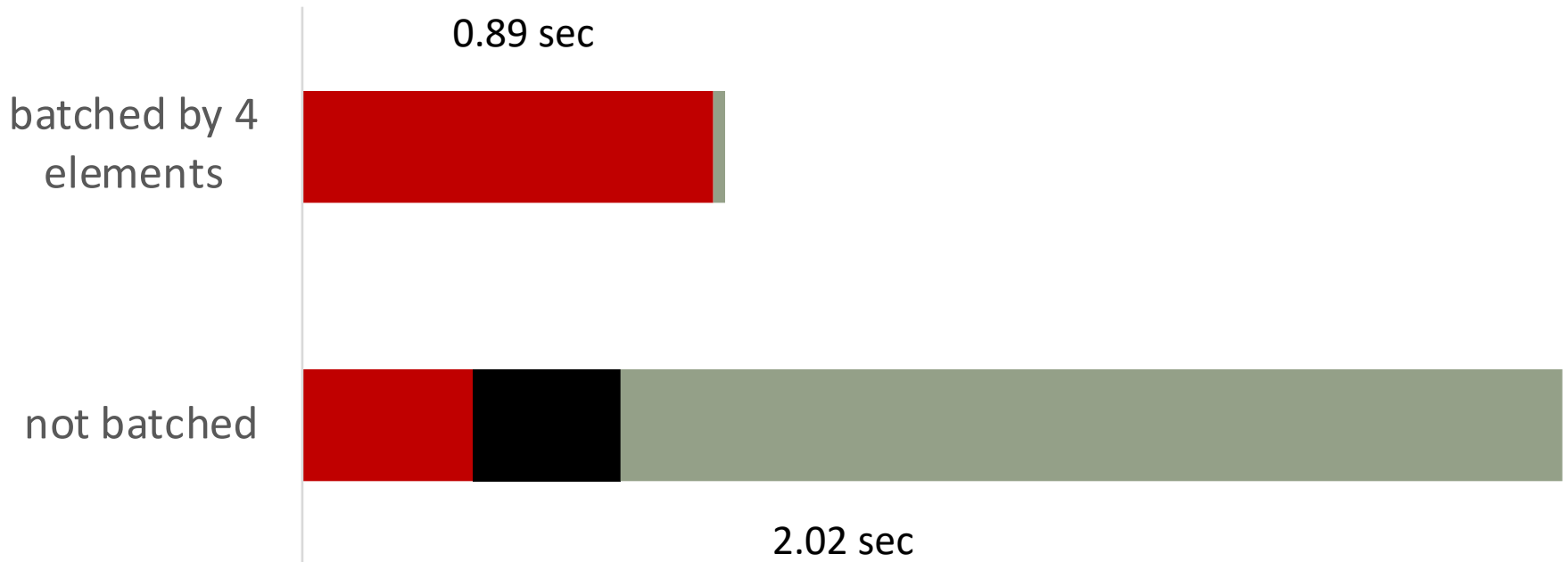batched by 4 elements
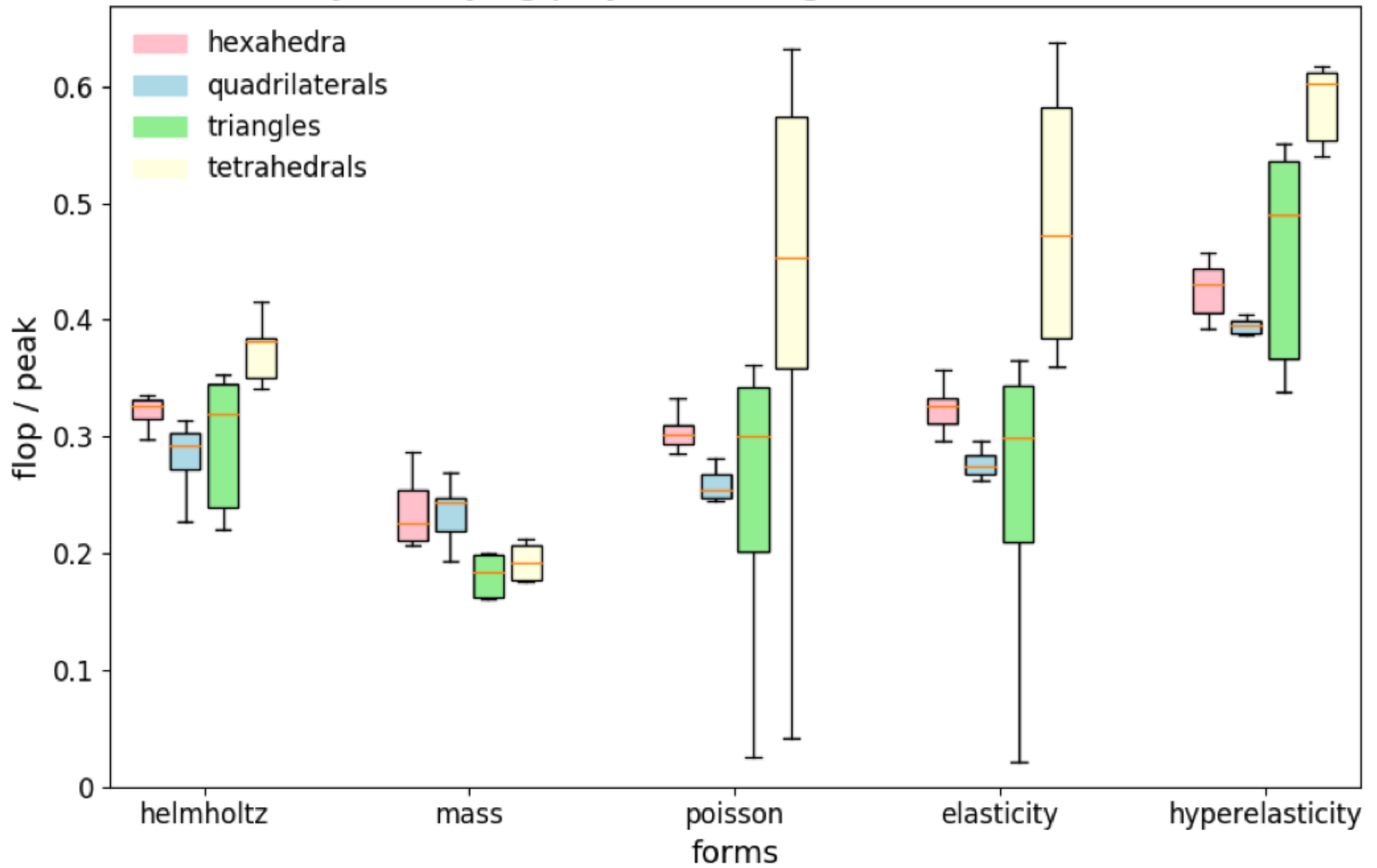
not batched

# Instruction counts by instruction types

■ AVX2 (4 doubles)    ■ SSE (2 doubles)    ■ scalar (1 double)

0.89 sec

batched by 4 elements

not batched

2.02 sec

batched by 4, varying polynomial degree for each form and mesh

# To be continued…

o We are building an abstraction layer of loops (via loo.py)

o Pathway to GPUs
   o …which requires a better performance model

o Try it out:
   o firedrake branch `tsfc2loopy`
   o tsfc branch `tsfc2loopy`
   o PyOP2 branch `tsfc_loopying`
   o `loopy branch opaque-types`

```
 1   static double const form_t19[3] = ...
 2   ...
 3
 4   void wrap_form00_cell_integral_otherwise(int const start, int const end, Mat const mat0,
 5     double const *__restrict__ dat0, int const *__restrict__ map0)
 6   {
 7     double t2[3 * 3] __attribute__ ((aligned (64)));
 8     ...
 9
10     for (int n = start; n <= -1 + end; ++n)
11     {
12       for (int i13 = 0; i13 <= 2; ++i13)
13         for (int i14 = 0; i14 <= 1; ++i14)
14           t3[2 * i13 + i14] = dat0[2 * map0[3 * n + i13] + i14];
15       for (int i1 = 0; i1 <= 2; ++i1)
16         for (int i4 = 0; i4 <= 2; ++i4)
17           t2[3 * i1 + i4] = 0.0;
18       ...
19       for (int form_j_0 = 0; form_j_0 <= 2; ++form_j_0)
20         for (int form_ip_0 = 0; form_ip_0 <= 2; ++form_ip_0)
21         {
22           form_t23 = form_t22[3 * form_ip_0 + form_j_0] * form_t9[form_ip_0];
23           for (int form_k_0 = 0; form_k_0 <= 2; ++form_k_0)
24             form_t24[3 * form_j_0 + form_k_0] += form_t22[3 * form_ip_0 + form_k_0] * form_t23;
25         }
26       ...
27       for (int form_k_2 = 0; form_k_2 <= 2; ++form_k_2)
28         for (int form_j_1 = 0; form_j_1 <= 2; ++form_j_1)
29           t2[3 * form_j_1 + form_k_2] += t2[3 * form_j_1 + form_k_2] + ...;
30
31       MatSetValuesBlockedLocal(mat0, 3, &(map0[3 * n]), 3, &(map0[3 * n]), &(t2[0]), ADD_VALUES);
32     }
33   }
```

CG1 Poisson bilinear form