

# Cross-Element Vectorisation in Firedrake

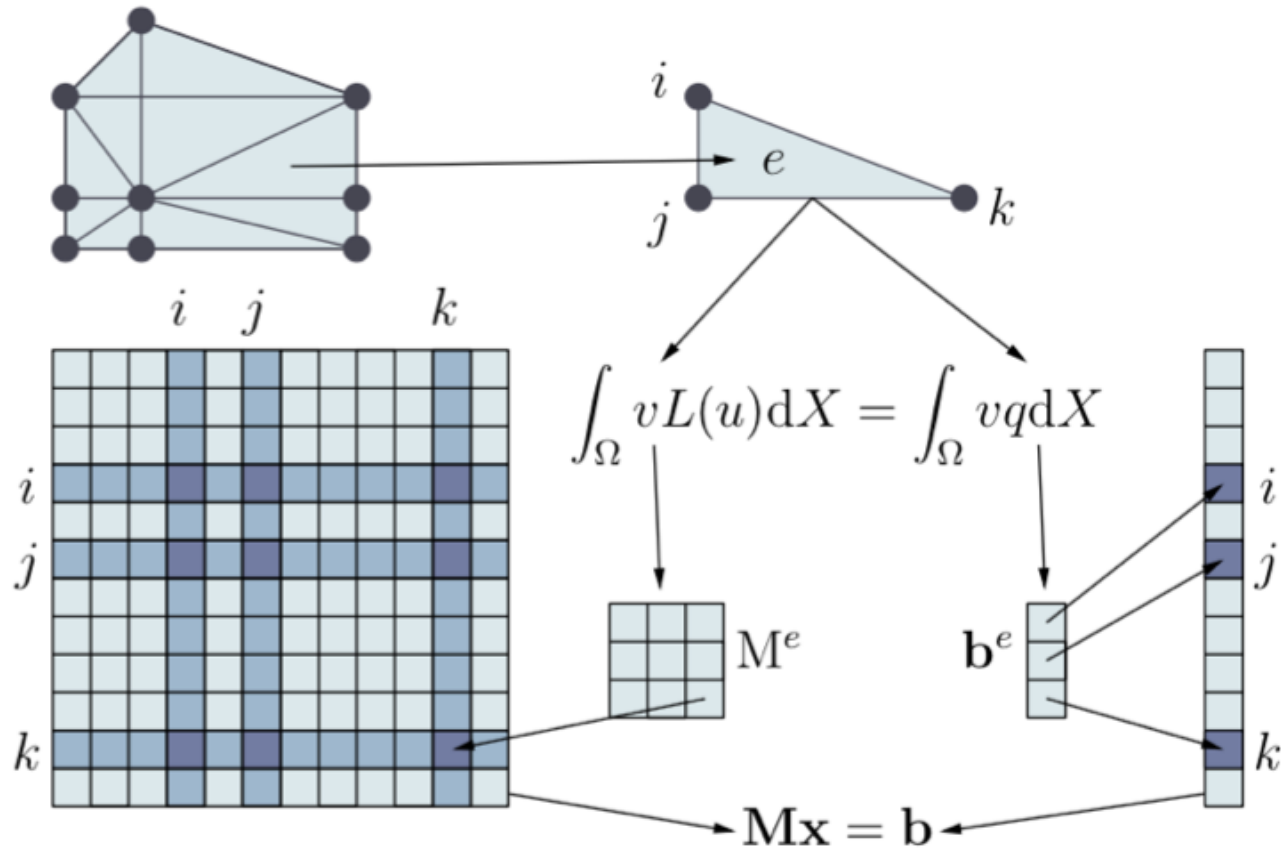
---

TJ Sun (ts2914@ic.ac.uk), Lawrence Mitchell, Kaushik Kulkarni, Andreas Klöckner  
David A Ham, Paul H J Kelly

April 2019

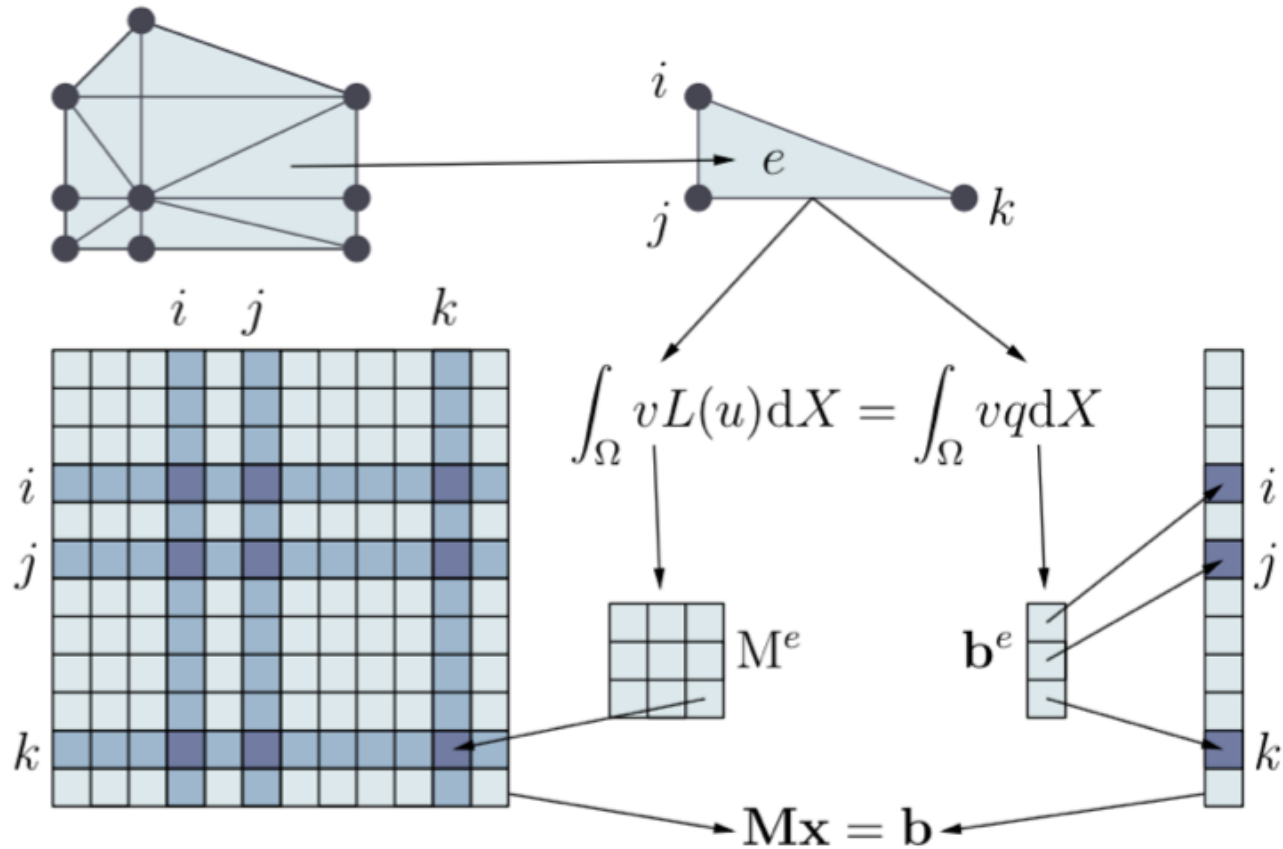
# Finite element method

(Computationally) finite element assembly  $\approx$  numerical integration



# Finite element method

(Computationally) finite element assembly  $\approx$  numerical integration

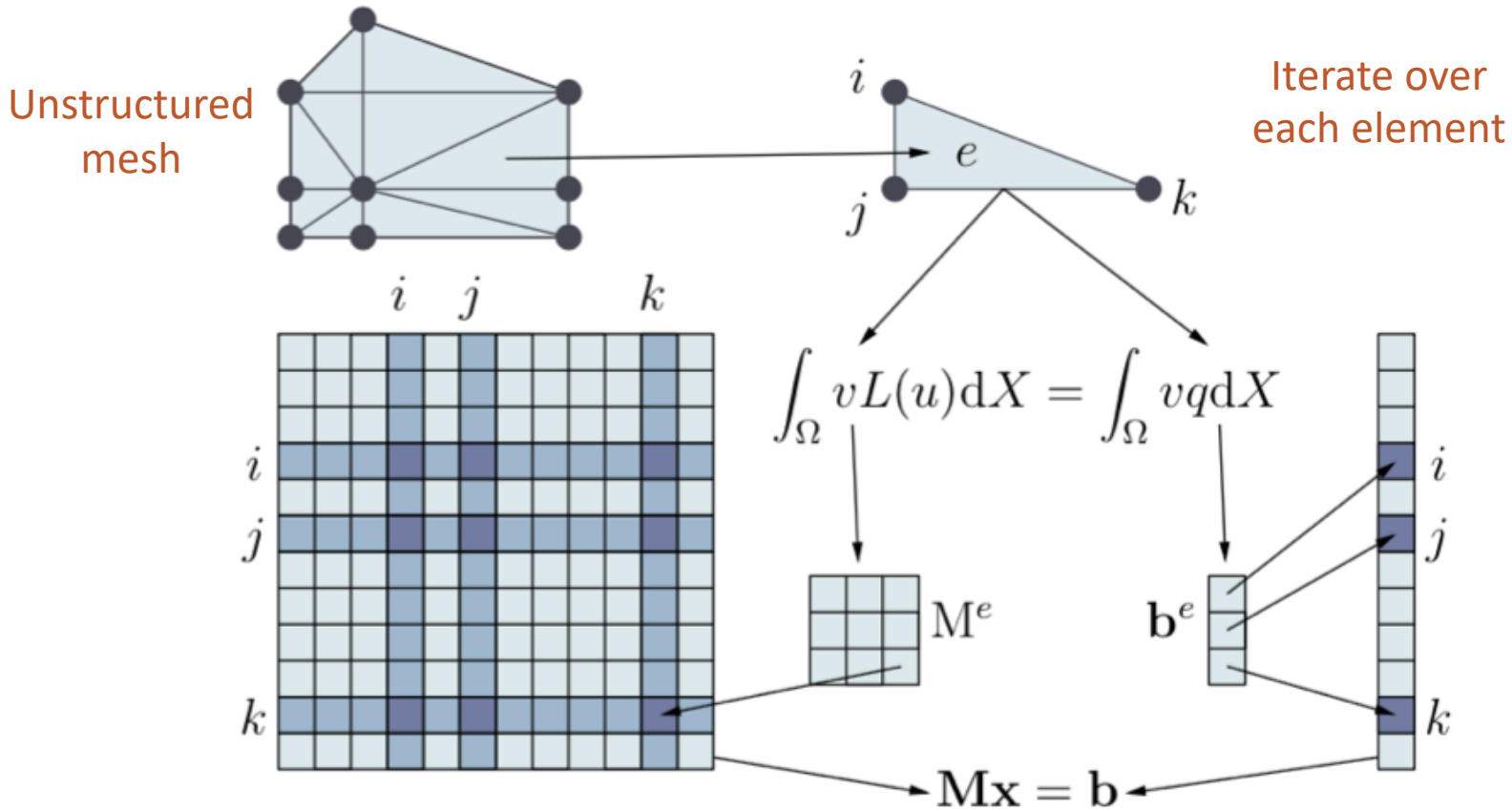


2-form assembly

1-form assembly

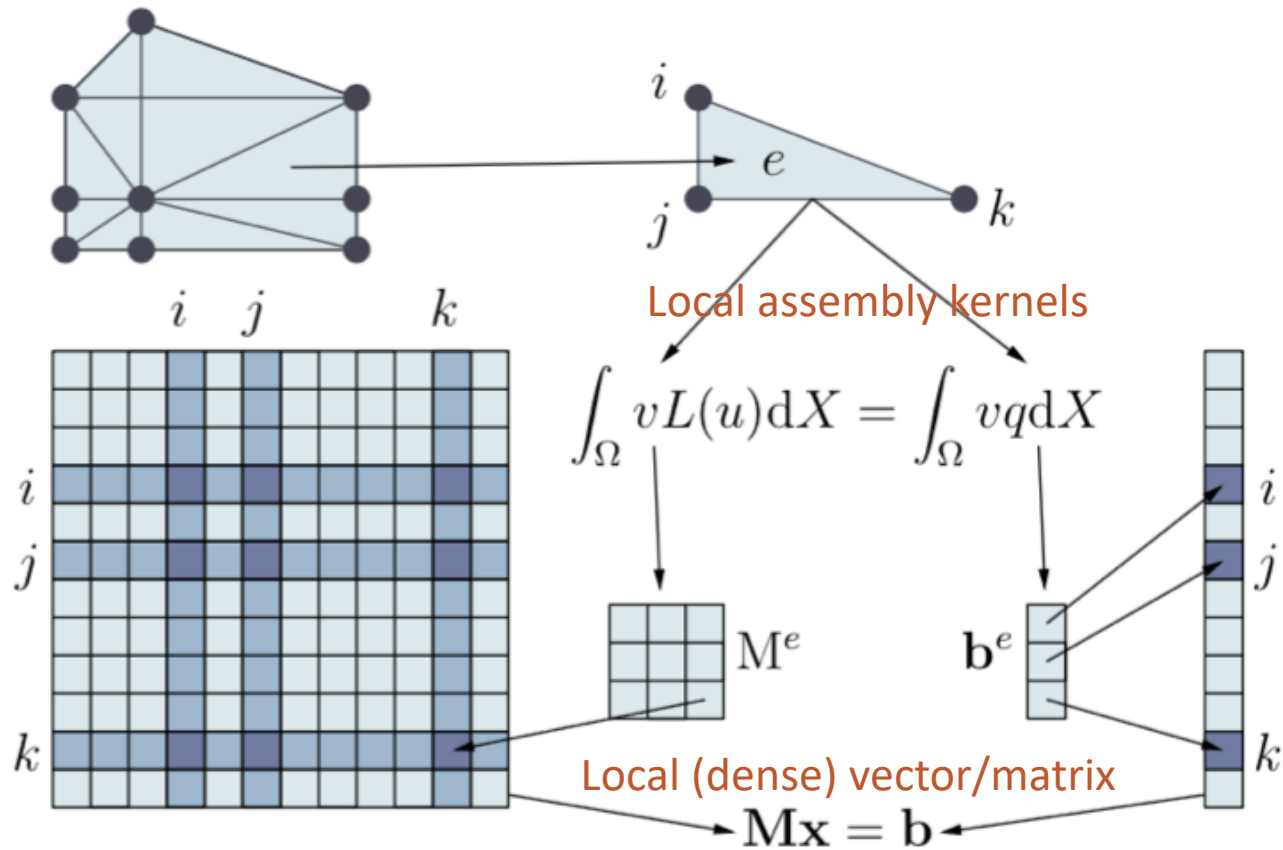
# Finite element method

(Computationally) finite element assembly  $\approx$  numerical integration



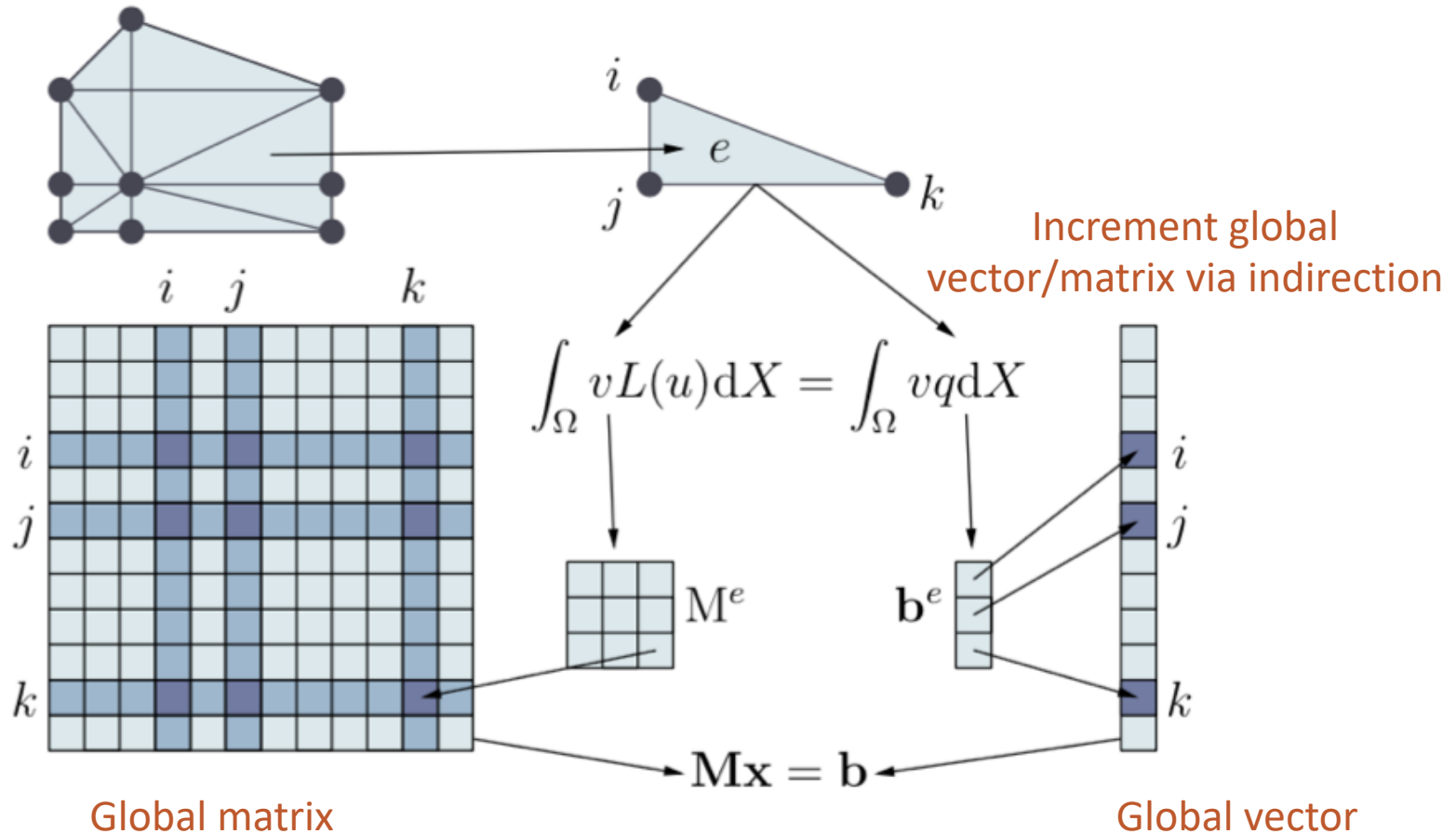
# Finite element method

(Computationally) finite element assembly  $\approx$  numerical integration



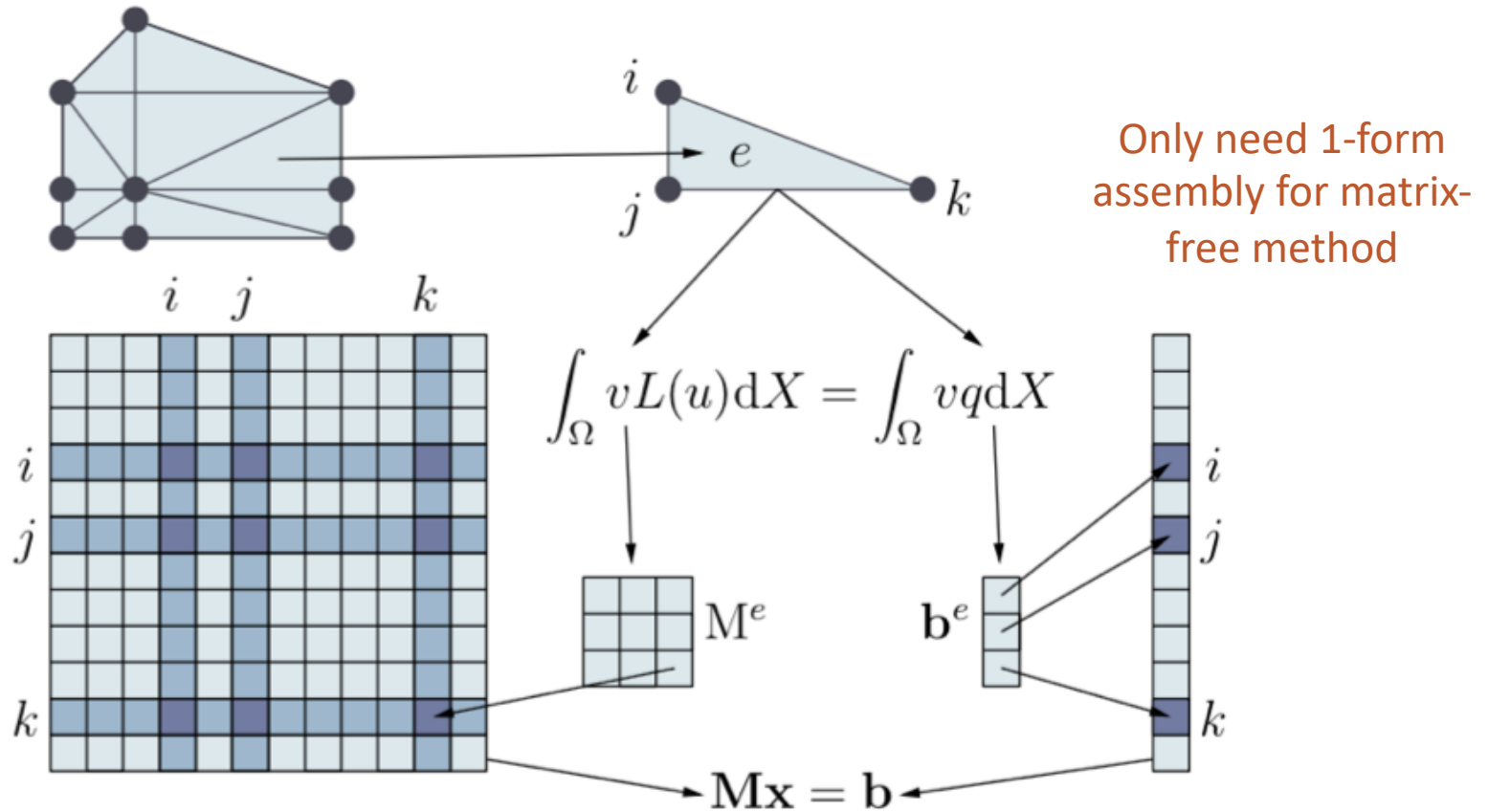
# Finite element method

(Computationally) finite element assembly  $\approx$  numerical integration



# Finite element method

(Computationally) finite element assembly  $\approx$  numerical integration



# Firedrake example

---

```
from firedrake import *

mesh = UnitSquareMesh(10, 10)
V = FunctionSpace(mesh, "CG", 1)
u = TrialFunction(V)
v = TestFunction(V)

# define the RHS
f = Function(V)
x, y = SpatialCoordinate(mesh)
f.interpolate((1+8*pi*pi)*cos(x*pi*2)*cos(y*pi*2))

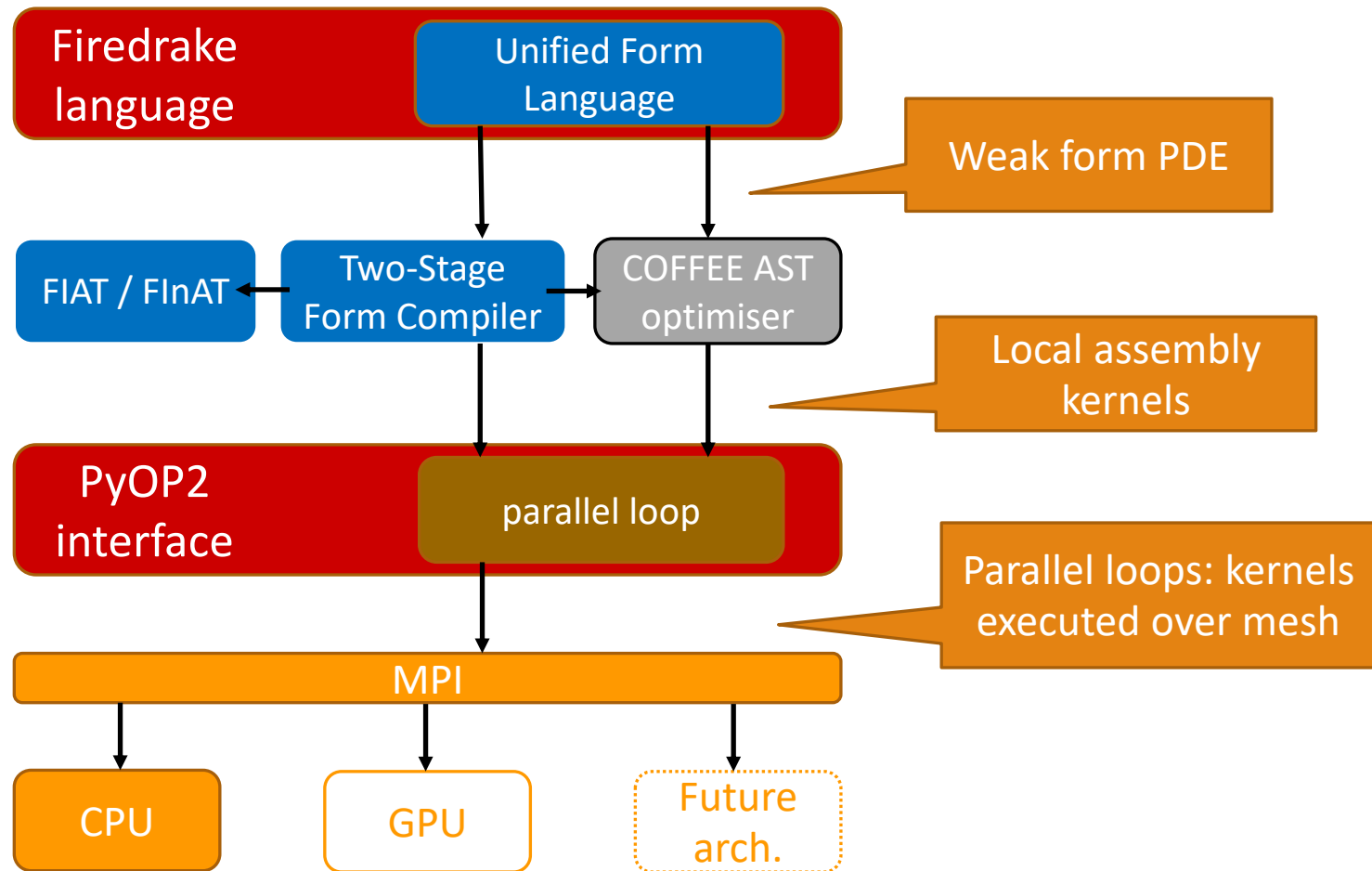
# define the PDE
a = (dot(grad(v), grad(u)) + u*v) * dx
L = f * v * dx

u = Function(V)
solve(a == L, u, solver_parameters={'ksp_type': 'cg'})

plot(u)
```

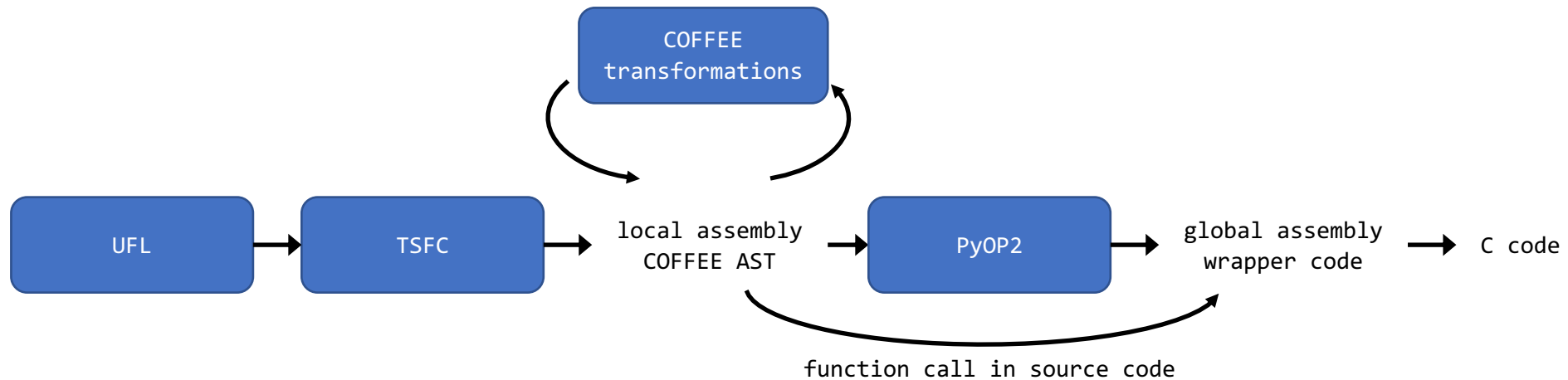


# Firedrake overview (simplified)



# “Intra-kernel” vectorisation of FE assembly

---



- Current default in Firedrake
- Vectorisation of **local** assembly kernel (AST) by COFFEE
- Challenges:
  - Loop trip counts can be small and/or not multiple of SIMD width
  - Usually requires alignment to cache boundary and stride 1 access
  - Operations outside of innermost loop not vectorized
  - Hard to automate for different PDEs, discretisations, meshes

# Cross-element vectorisation

---

- Vectorisation of **global** assembly
- Vector-expand the local assembly kernel to act on N elements together, with  $N = \text{SIMD width}$
- All operations can be vectorised
- We can always do this systematically
- But this requires:
  - Intermediate Representation of the whole global assembly
  - Robust loop transformations

# Cross-element vectorisation

---

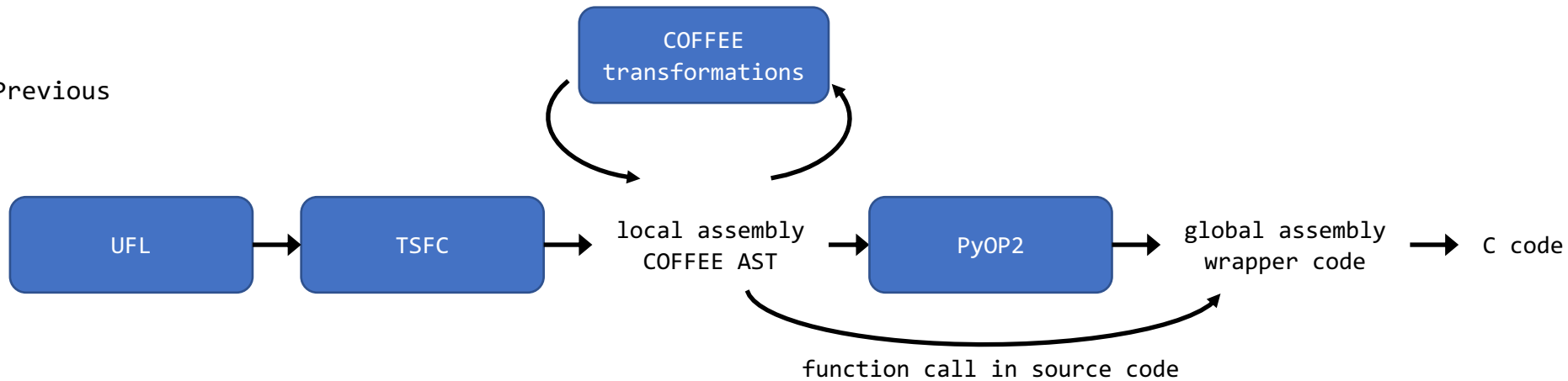
- Vectorisation of **global** assembly
- Vector-expand the local assembly kernel to act on N elements together, with  $N = \text{SIMD width}$
- All operations can be vectorised
- We can always do this systematically
- But this requires:
  - Intermediate Representation of the whole global assembly
  - Robust loop transformations

Introducing **Loopy** as our new code generation backend

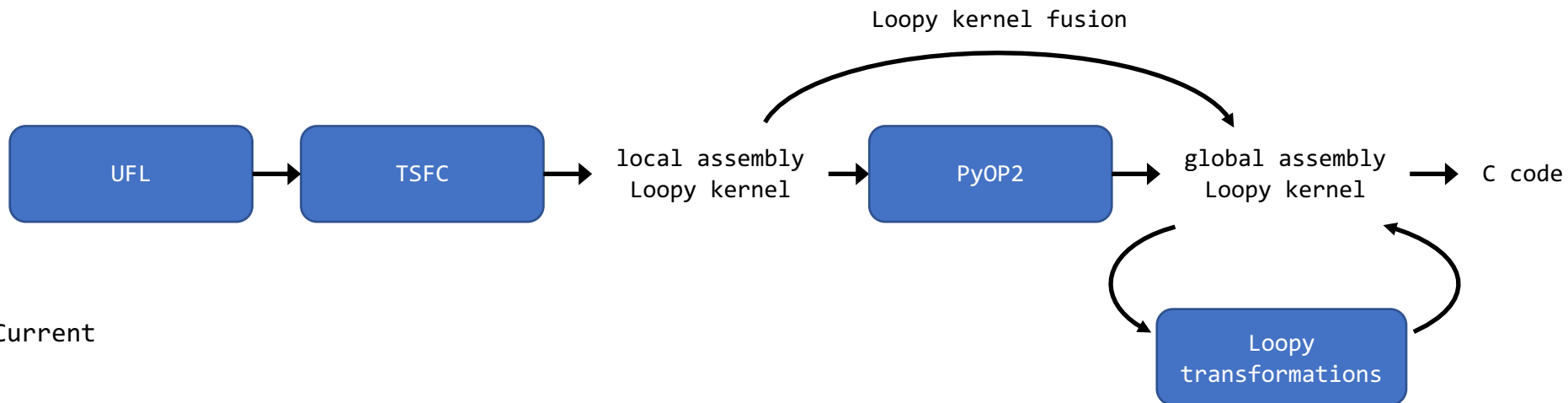
- Python package by Andreas Klöckner et al. (UIUC)
- Based on polyhedral model of loops
- $\approx$  model of loops + transformations
- “Do what I say”, not a black-box loop optimiser
- Support multiple backends from same IR: CPU, ISPC, OpenCL, CUDA

# Code generation in Firedrake

Previous



Current



```

1 // ... Constant array declarations ... //
2
3 void wrap_helmholtz(int const start, int const end, double *__restrict__ dat0, double const
  *__restrict__ dat1, double const *__restrict__ dat2, int const *__restrict__ map0)
4 {
5     double t2[3];
6     double t3[3 * 2];
7     double t4[3];
8
9     for (int n = start; n <= -1 + end; ++n)
10    {
11        for (int i5 = 0; i5 <= 2; ++i5)
12        {
13            t4[i5] = dat2[map0[3 * n + i5]];
14            for (int i6 = 0; i6 <= 1; ++i6)
15                t3[2 * i5 + i6] = dat1[2 * map0[3 * n + i5] + i6];
16        }
17        for (int i1 = 0; i1 <= 2; ++i1)
18            t2[i1] = 0.0;
19
20
21        double form_t11 = 0.0;
22        double form_t1 = -1.0 * t3[0];
23        // ... More similar instructions ... //
24
25        for (int i15 = 0; i15 <= 2; ++i15)
26            dat0[map0[3 * n + i15]] += t2[i15];
27    }
28 }

```

```

1 // ... Constant array declarations ... //
2
3 void wrap_helmholtz(int const start, int const end, double *__restrict__ dat0, double const
  *__restrict__ dat1, double const *__restrict__ dat2, int const *__restrict__ map0)
4 {
5     double t3[3];
6     double t2[3];
7     double t4[3];
8
9     for (int n = start; n <= -1 + end; ++n)
10    {
11        for (int i5 = 0; i5 <= 2; ++i5)
12        {
13            t4[i5] = dat2[map0[3 * n + i5]];
14            for (int i6 = 0; i6 <= 1; ++i6)
15            {
16                t3[2 * i5 + i6] = dat1[2 * map0[3 * n + i5] + i6];
17            }
18            for (int i1 = 0; i1 <= 2; ++i1)
19            {
20                t2[i1] = 0.0;
21            }
22
23            double form_t11 = 0.0;
24            double form_t1 = -1.0 * t3[0];
25            // ... More similar instructions ... //
26
27            for (int i15 = 0; i15 <= 2; ++i15)
28            {
29                dat0[map0[3 * n + i15]] += t2[i15];
30            }
31        }
32    }
33 }

```

Outer loop over all elements in the mesh

Data gathering

Inlined local kernel

Data scattering

# The Loopy representation

---

## KERNEL:

helmholtz

## ARGUMENTS:

```
start: type: int32
end: type: int32
dat0: type: float64, shape: (None)
dat1: type: float64, shape: (None, 2)
dat2: type: float64, shape: (None)
map0: type: int32, shape: (None, 3)
```

## DOMAINS:

```
[end, start] -> { [n] : start <= n < end }
{ [i2] : 0 <= i2 <= 2 }
// ... More domains ... //
```

## INAME\_IMPLEMENTATION\_TAGS:

None

## TEMPORARIES:

```
t4: type: float64, shape: (3), dim_tags: (stride:1)
// ... More temporaries ... //
```

## INSTRUCTIONS:

```
for n, i2
    t4[i2] = dat2[map0[n, i2]]
    for i3
        t1[i2, i3] = dat1[map0[n, i2], i3]
    end i2, i3
    for i0
        t0[i0] = 0
    end i0
    // ... Inlined local assembly kernel ... //
    form_t11 = 0.0
    // ... More instructions ... //
    for i15
        dat0[map0[n, i15]] += t0[0, i15]
    end n, i15
```



# The Loopy representation

## KERNEL:

helmholtz

## ARGUMENTS:

start: type: int32

end: type: int32

map0: type: int32, shape: (None)

map1: type: int32, shape: (None, 2)

map2: type: int32, shape: (None)

map0: type: int32, shape: (None, 3)

## DOMAINS:

[end, start] -> { [n] : start <= n < end }

{ [i2] : 0 <= i2 <= 2 }

// ... More domains ... //

## INAME\_IMPLEMENTATION\_TAGS:

None

## TEMPORARIES:

t4: type: float64, shape: (3)

// ... More temporaries ... //

## INSTRUCTIONS:

for n, i2

t4[i2] = dat2[map0[n, i2]]

for i3

t1[i2, i3] = dat1[map0[n, i2], i3]

end i2, i3

for i0

t0[i0] = 0

assembly kernel ... //

// ... More instructions ... //

for i15

dat0[map0[n, i15]] += t0[0, i15]

end n, i15

Loop indices are called "inames"

Affine inequalities describe the loop bounds

Tags governs code generation behaviour

# Cross-element vectorisation by applying Loopy transformations

---

1. Split the outer loop  $n$  into  $n_{\text{outer}}$  and  $n_{\text{simd}}$ 
  - With  $n_{\text{simd}}$  extent = SIMD width
2. Tag  $n_{\text{simd}}$  with implementation tag “SIMD”

```

1 // ... Constant array declarations ... //
2
3 void wrap_helmholtz(int const start, int const end, double *__restrict__ dat0, double const *
  __restrict__ dat1, double const *__restrict__ dat2, int const *__restrict__ map0)
4 {
5     double t2[3 * 4] __attribute__((aligned (64)));
6     // ... More temporary array declarations ... //
7     for (int n_out = (start / 4); n_out <= ((end-4) / 4); ++n_out) {
8         for (int i5 = 0; i5 <= 2; ++i5) {
9             for (int i6 = 0; i6 <= 1; ++i6) {
10                 #pragma omp simd
11                 for (int n_simd = 0; n_simd <= 3; ++n_simd)
12                     t3[n_simd+8*i5+4*i6]=dat1[2*map0[12*n_out+3*n_simd+i5]+i6];
13             }
14             #pragma omp simd
15             for (int n_simd = 0; n_simd <= 3; ++n_simd)
16                 t4[n_simd+4*i5] = dat2[map0[12*n_out+3*n_simd+i5]];
17         }
18         for (int i1 = 0; i1 <= 2; ++i1) {
19             #pragma omp simd
20             for (int n_simd = 0; n_simd <= 3; ++n_simd)
21                 t2[n_simd + 4 * i1] = 0.0;
22         }
23         #pragma omp simd
24         for (int n_simd = 0; n_simd <= 3; ++n_simd) {
25             form_t11[n_simd] = 0.0;
26             form_t1[n_simd] = -1.0 * t3[n_simd];
27             // ... More similar instructions ... //
28         }
29         for (int i15 = 0; i15 <= 2; ++i15)
30             for (int n_batch = 0; n_batch <= 3; ++n_batch)
31                 dat0[map0[12*n_out+3*n_batch+i15]] += t2[n_batch+4*i15];
32     }
33 }

```

Residual assembly of Helmholtz operator on triangle mesh, CG1 space, batched by 4

# Temporaries vector-expanded, with n\_simd dimension fastest moving

```
1 // Constant array declarations //
2
3 __restrict__ dat1, double const *__restrict__ dat2, int const *__restrict__ map0)
4 {
5     double t2[3 * 4] __attribute__((aligned(64)));
6     // ... More temporary array declarations ... //
7     for (int n_out = (start / 4); n_out <= ((end-4) / 4); ++n_out) {
8         for (int
9             for (int
10                 #pragma omp simd
11                 for (int n_simd = 0; n_simd <= 3; ++n_simd)
12                     t3[n_simd+8*i5+4*i6]=dat1[2*map0[12*n_out+3*n_simd+i5]+i6];
13             }
14             #pragma omp simd
15             for (int n_simd = 0; n_simd <= 3; ++n_simd)
16                 t4[n_simd+4*i5] = dat2[map0[12*n_out+3*n_simd+i5]];
17         }
18         for (int i1 = 0; i1 <= 2; ++i1) {
19             #pragma omp simd
20             for (int n_simd = 0; n_simd <= 3; ++n_simd)
21                 t2[n_simd + 4 * i1] = 0.0;
22         }
23         #pragma omp simd
24         for (int n_simd = 0; n_simd <= 3; ++n_simd) {
25             form_t11[n_simd] = 0.0;
26             form_t1[n_simd] = -1.0 * t3[n_simd];
27             // ... More ...
28         }
29         for (int i1 = 0; i1 <= 2; ++i1,
30             for (int n_batch = 0; n_batch <= 3; ++n_batch)
31                 dat0[map0[12*n_out+3*n_batch+i15]] += t2[n_batch+4*i15];
32     }
33 }
```

n\_simd loops pushed to innermost, decorated by pragma

Write-back is serialised, due to potential race condition

Residual assembly of Helmholtz operator on triangle mesh, CG1 space, batched by 4

# Compiler vector extensions

---

- A more direct way to inform compiler to generate SIMD instructions than using OpenMP pragma
- For example
  - `typedef double double4 __attribute__((vector_size(32)));`
- We need to add new features to Loopy
  - New code generation target with vector data types
  - New iname tag

```

1  typedef double double4 __attribute__((vector_size (32)));
2  typedef int int4 __attribute__((vector_size (16)));
3
4  // ... Constant array declarations ... //
5
6  void wrap_form0_cell_integral_otherwise(int const start, int const end, double *__restrict__ dat0,
    double const *__restrict__ dat1, double const *__restrict__ dat2, int const *__restrict__ map0)
7  {
8      double4 form_t1 __attribute__((aligned (64)));
9      // ... Temporary array declarations ... //
10
11     for (int n_outer = (start / 4); n_outer <= ((-4 + end) / 4); ++n_outer) {
12         for (int i5 = 0; i5 <= 2; ++i5) {
13             for (int i6 = 0; i6 <= 1; ++i6) {
14                 #pragma omp simd
15                 for (int n_simd = 0; n_simd <= 3; ++n_simd)
16                     t3[(2 * i5 + i6)][n_simd] = dat1[2 * map0[12 * n_outer + 3 * n_simd + i5] + i6];
17             }
18             #pragma omp simd
19             for (int n_simd = 0; n_simd <= 3; ++n_simd)
20                 t4[i5][n_simd] = dat2[map0[12 * n_outer + 3 * n_simd + i5]];
21         }
22         for (int i1 = 0; i1 <= 2; ++i1)
23             t2[i1] = _zeros_double4;
24
25         form_t11 = _zeros_double4;
26         form_t1 = -1.0 * t3[0];
27         // ... More similar instructions ... //
28
29         for (int i15 = 0; i15 <= 2; ++i15)
30             for (int n_batch = 0; n_batch <= 3; ++n_batch)
31                 dat0[map0[12 * n_outer + 3 * n_batch + i15]] += t2[i15][n_batch];
32     }
33 }

```

Residual assembly of Helmholtz operator on triangle mesh, CG1 space, vector extensions

```

1 typedef double double4 __attribute__((vector_size (32)));
2 typedef int int4 __attribute__((vector_size (16)));
3
4 // ... Constant array declarations ...
5
6 void wrap_form0_cell_integral_generic(const start, const end, double *__restrict__ dat0,
7   double const *__restrict__ dat1, double const *__restrict__ dat2, int const *__restrict__ map0)
8 {
9   double4 form_t1 __attribute__((aligned (64)));
10   // ... Temporary array declarations ... //
11   for (int n_outer = (start / 4); n_outer <= ((-4 + end) / 4); ++n_outer) {
12     for (int i5 = 0; i5 <= 2; ++i5) {
13       for (int i6 = 0; i6 <= 1; ++i6) {
14         #pragma omp simd
15         for (int n_simd = 0; n_simd <= 3; ++n_simd)
16           t3[(2 * i5 + i6)][n_simd] = dat1[2 * map0[12 * n_outer + 3 * n_simd + i5] + i6];
17       }
18       #pragma omp simd
19       for (int n_simd = 0; n_simd <= 3; ++n_simd)
20         t4[i5][n_simd] = dat2[map0[12 * n_outer + 3 * n_simd + i5]];
21     }
22     for (int i1 = 0; i1 <= 2; ++i1)
23       t2[i1] = _zeros_double4;
24
25     form_t11 = _zeros_double4;
26     form_t1 = -1.0 * t3[0];
27     // ... More similar instructions ... //
28
29     for (int i15 = 0; i15 <= 2; ++i15)
30       for (int n_batch = 0; n_batch <= 3; ++n_batch)
31         dat0[map0[12 * n_outer + 3 * n_batch + i15]] += t2[i15][n_batch];
32   }
33 }

```

**Vector data types**

**Vector instructions**

Residual assembly of Helmholtz operator on triangle mesh, CG1 space, vector extensions

# Experimental setup

---

- Hardware

	Haswell Xeon E5-2640 v3	Skylake Xeon Gold 6148
Base frequency	2.6 GHz	2.4 GHz
Physical cores	8	20
SIMD instruction set	AVX2	AVX512
doubles per SIMD vector	4	8
Cross-element vectorization batch size	4	8
FMA <sup>4</sup> units per core	2	2
FMA instruction issue per cycle	2	2
Peak performance (double-precision) <sup>5</sup>	332.8 GFLOP/s	1536.0 GFLOP/s
LINPACK performance (double-precision) <sup>6</sup>	262.5 GFLOP/s	976.7 GFLOP/s
Memory bandwidth <sup>7</sup>	38.5 GB/s	81.0 GB/s

- We measure all cores in a single node and report % of peak performance
- 4 meshes: triangles, quadrilaterals, tetrahedra, hexahedra
- 5 operators
- 3 compilers: GCC, ICC, Clang



# Check that SIMD instructions are generated

---

- Flop contributions by instruction types (Helmholtz on hexahedra, degree 4, ICC)

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

Only 0.4% of  
flops not  
vectorized

batched by 4  
elements

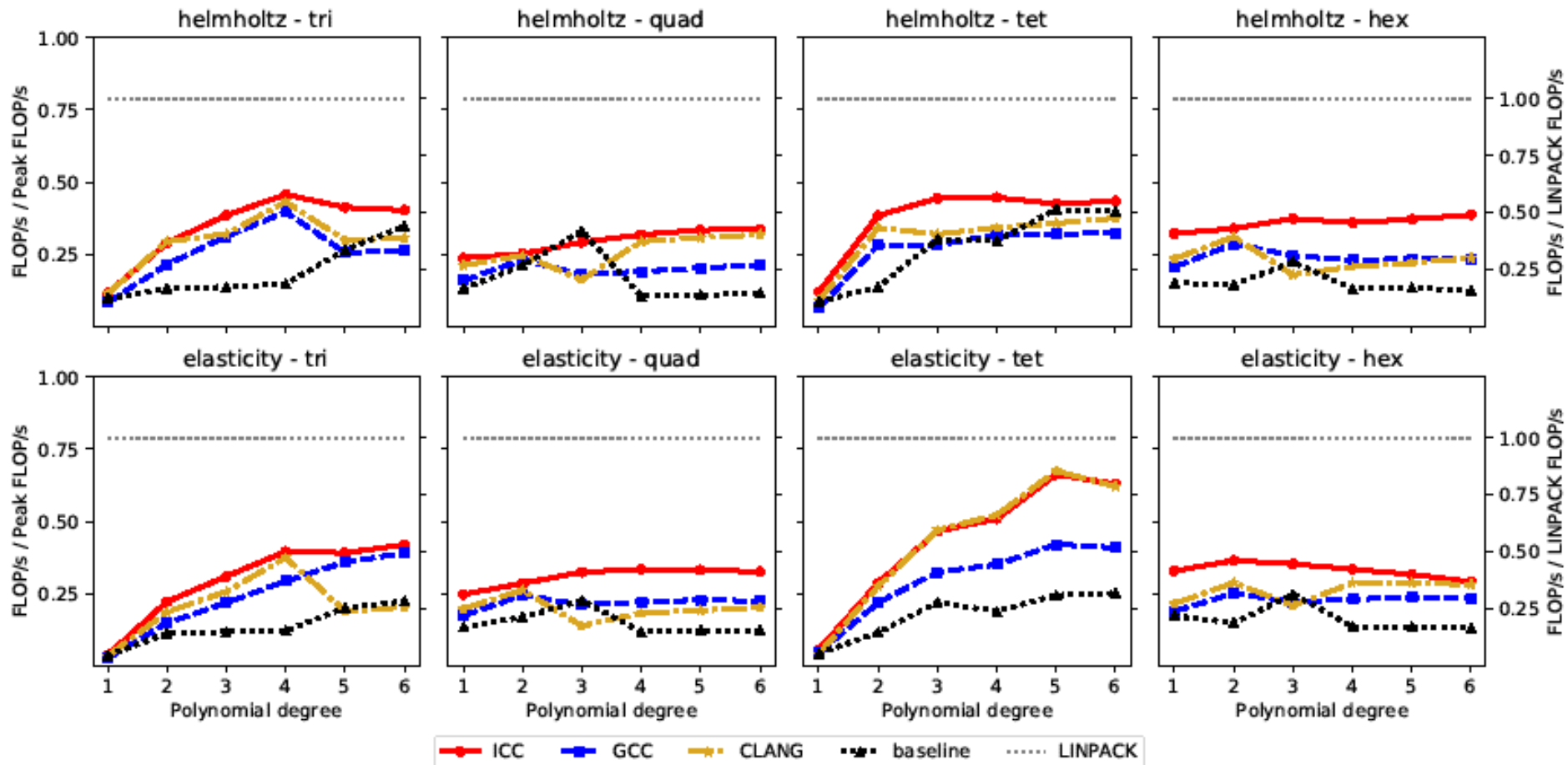


not batched



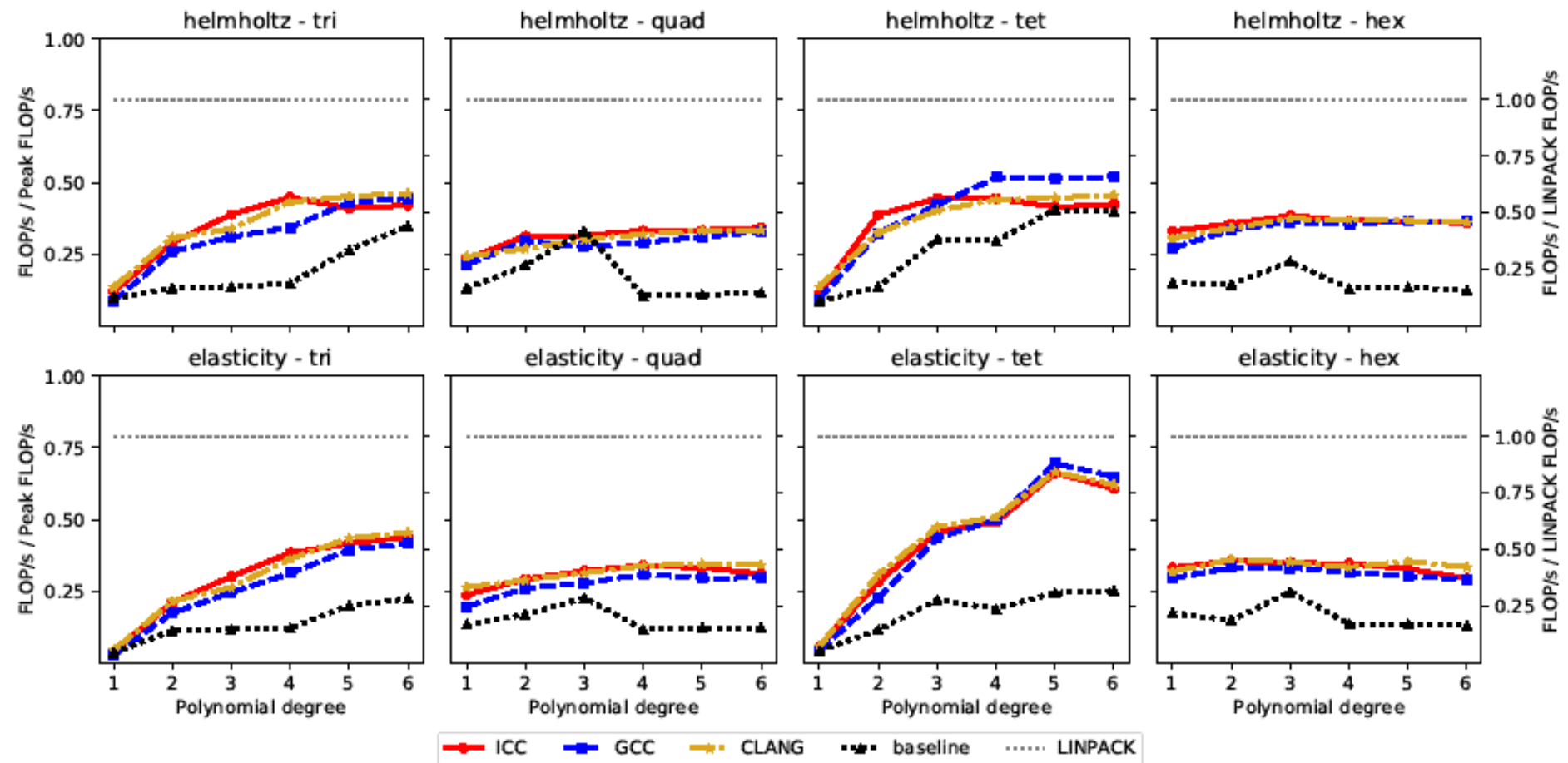
# Haswell, OpenMP pragma

- ICC is consistent across different polynomial degrees



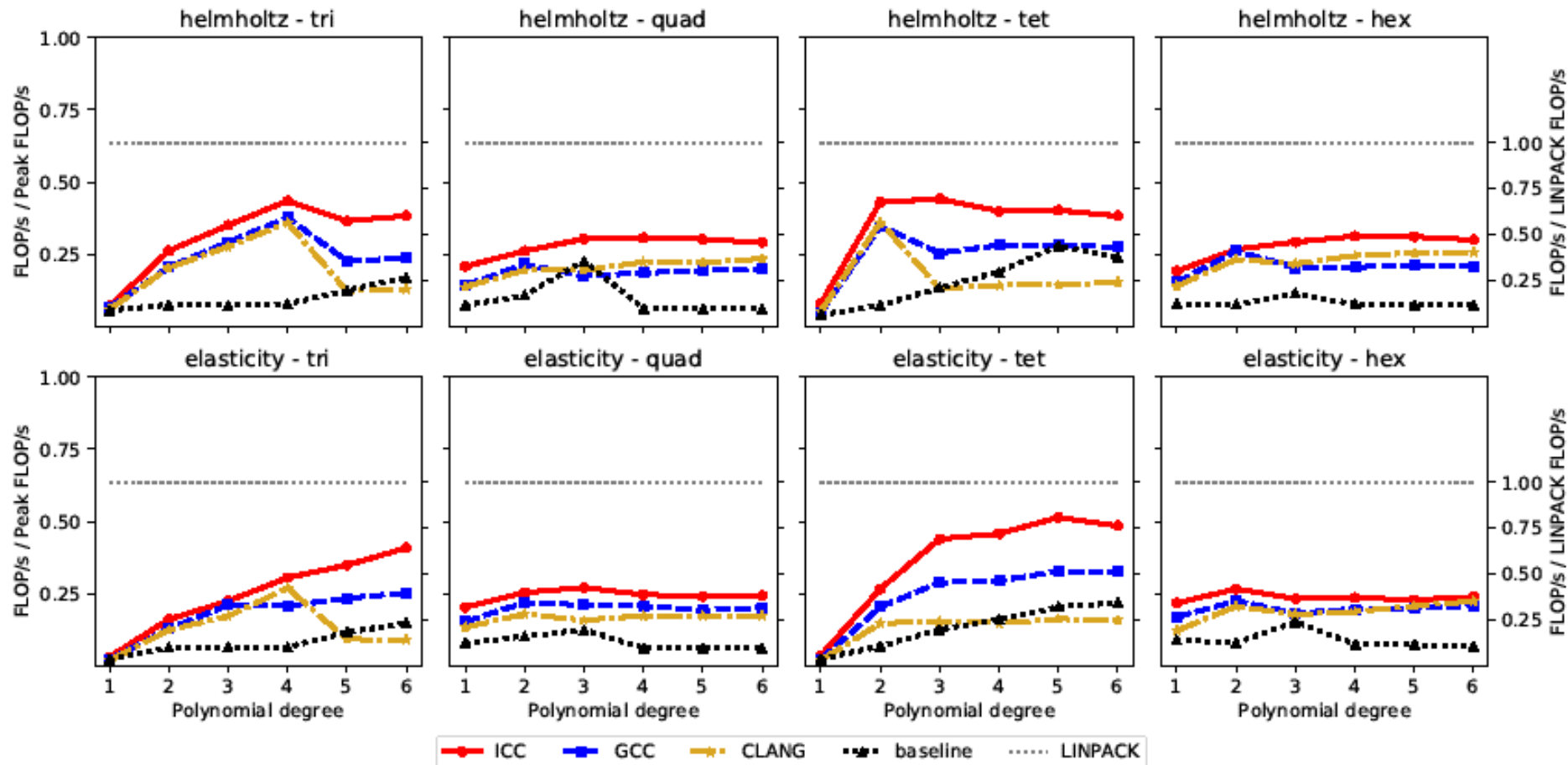
# Haswell, vector extensions

- Vector extensions give portable performance



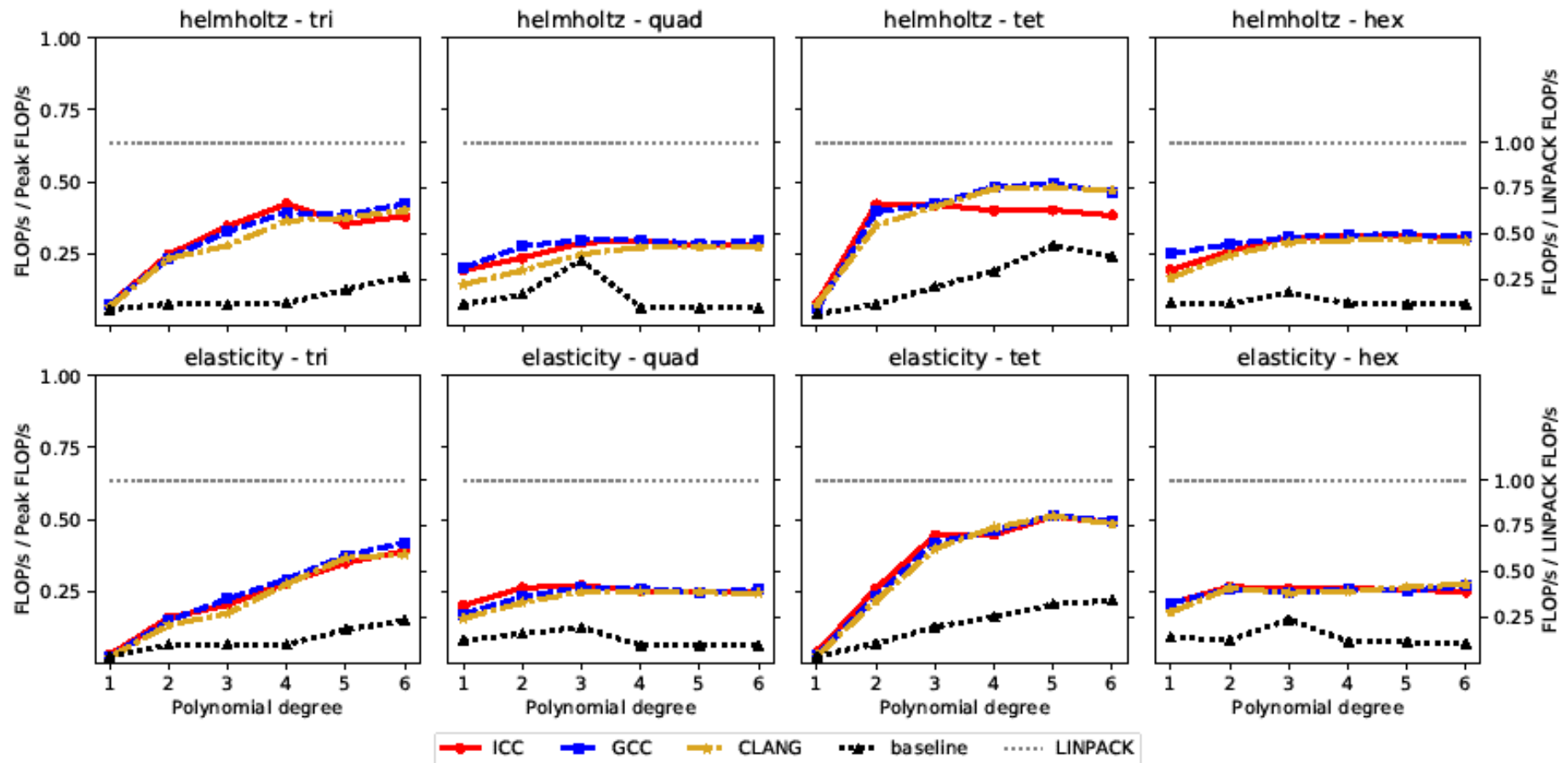
# Skylake, OpenMP pragma

- Similar observations as Haswell

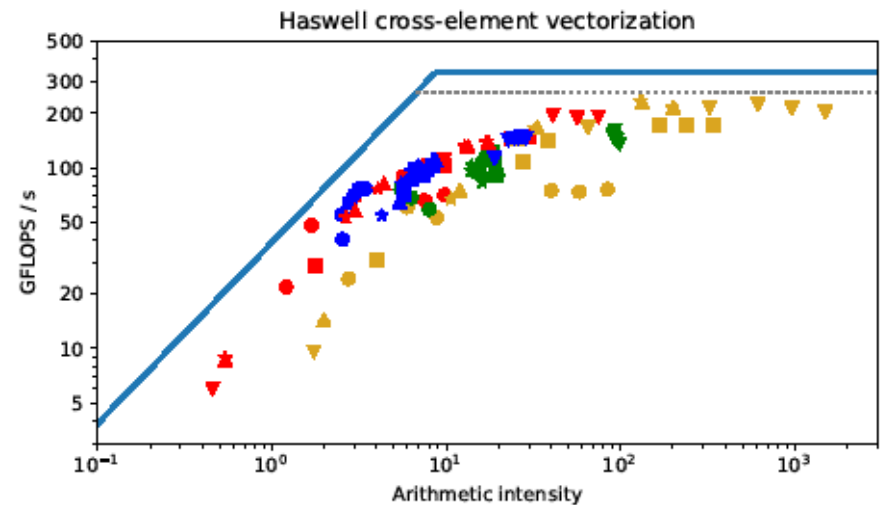
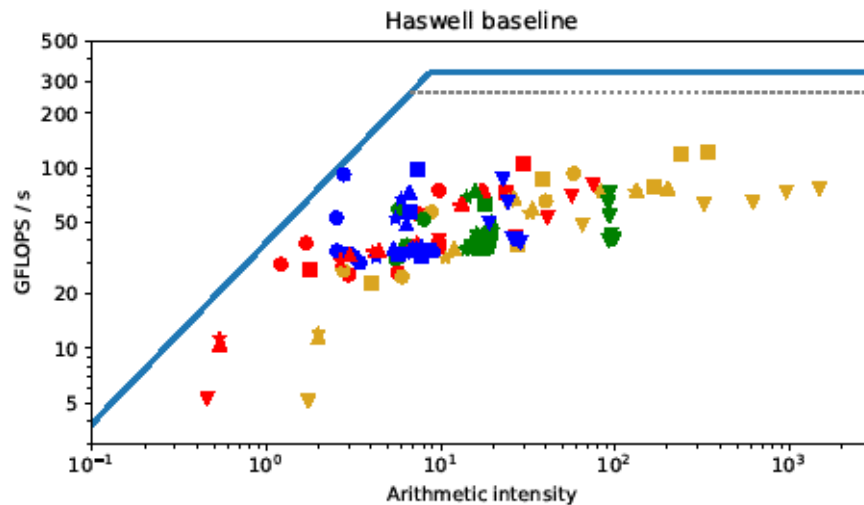


# Skylake, vector extensions

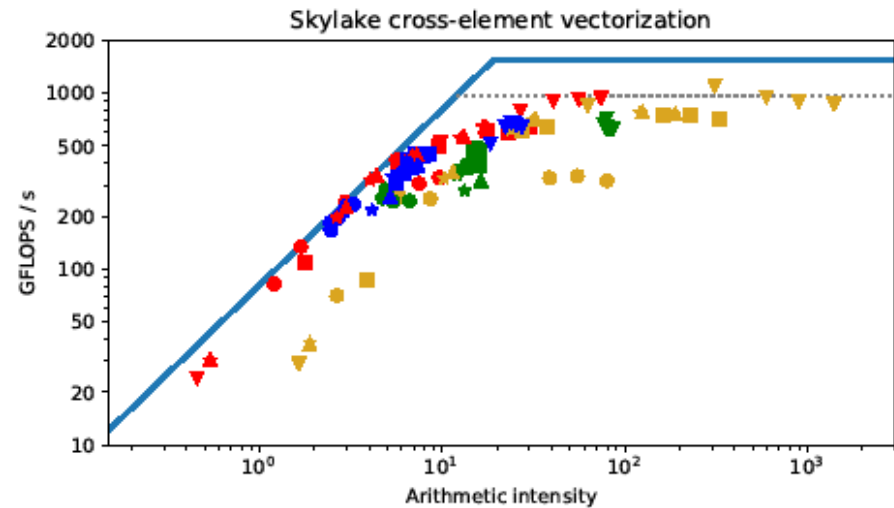
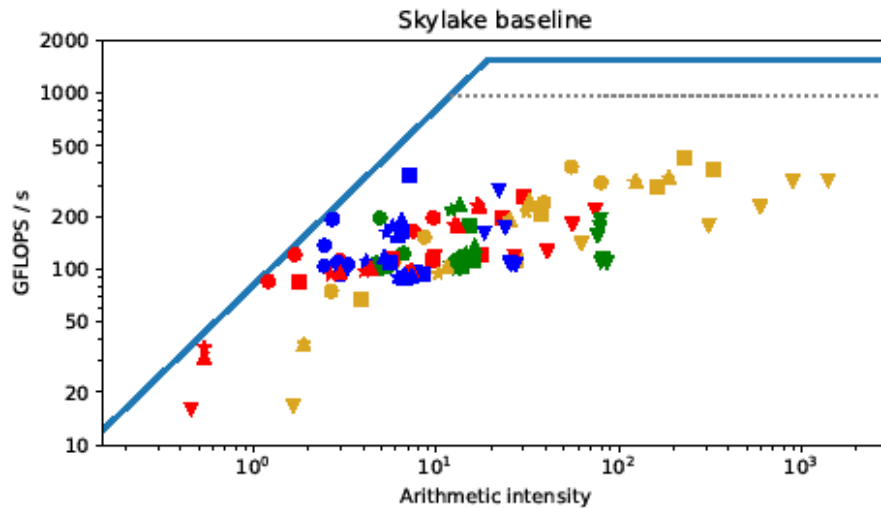
- Similar observations as Haswell



# Haswell roofline



# Skylake roofline



- |               |                    |                    |                     |                          |
|---------------|--------------------|--------------------|---------------------|--------------------------|
| ● mass - tri  | ■ helmholtz - tri  | ★ laplacian - tri  | ▲ elasticity - tri  | ▼ hyperelasticity - tri  |
| ● mass - quad | ■ helmholtz - quad | ★ laplacian - quad | ▲ elasticity - quad | ▼ hyperelasticity - quad |
| ● mass - tet  | ■ helmholtz - tet  | ★ laplacian - tet  | ▲ elasticity - tet  | ▼ hyperelasticity - tet  |
| ● mass - hex  | ■ helmholtz - hex  | ★ laplacian - hex  | ▲ elasticity - hex  | ▼ hyperelasticity - hex  |

..... LINPACK

# Conclusion

---

- Cross-element vectorisation consistently achieves >30% peak performance for residual assembly
- New abstraction for global assembly enables automated and robust transformation and code generation, powered by Loopy
- More to come...
  - Better support for bilinear form assembly (currently need to pack into individual element matrices and call PETSc function MatSetValues)
  - Pathway to GPUs (explore more vectorisation strategies other than “1 element per thread”)
- Preprint on arXiv: A study of vectorization for matrix-free finite element methods
- [www.firedrakeproject.org](http://www.firedrakeproject.org) for docs, demos, notebooks



*Firedrake*