

Towards cross-domain domain-specific compiler architecture

Paul Kelly
Group Leader, Software Performance Optimisation
Department of Computing
Imperial College London

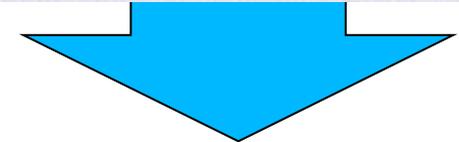
Joint work with David Ham (Imperial Maths), Lawrence Mitchell (University of Durham)
George Bisbas, Edward Stow, Fabio Luporini (Devito Codes Ltd), Florian Rathgeber (now with Google), Doru Bercea (now with IBM Research), Michael Lange (now with ECMWF), Andrew McRae (now at University of Oxford), Graham Markall (now at NVIDIA), Tianjiao Sun (now at Cerebras), Thomas Gibson (NCSA Illinois), Kaushik Kulkarni (UIUC), Andreas Klockner (UIUC), Tobias Grosser, Michel Steuwer (University of Edinburgh), Amrey Krause, Nick Brown (EPCC)
And many others....

Who am I and what do I do?

- I have worked on general-purpose compilers
 - Notably pointer analysis
 - adopted into GCC
 - (actually the work of my PhD student David Pearce)
- But the benefits were incremental
- Meanwhile I engaged with applications specialists
 - Who know they have major performance optimisation opportunities
 - So I got interested in automating domain-specific optimisations

This talk is about domain-specific languages

- So we can deliver domain-specific **optimisations**
- So we collect and automate all the performance techniques that are known for a **family of problems**
- If we get it right.... we get
 - **Productivity** – by generating low-level code from a high-level specification
 - **Performance** – by automating optimisations
 - **Performance portability** – with multiple back-ends



Have your cake and eat it too

- I aim to show you that you **can** simultaneously
 - **raise the level at** which programmers can reason about code,
 - **provide the compiler with a model of the computation that enables it to generate faster code than you could reasonably write by hand**



The work of my research group

- Vectorisation, parametric polyhedral tiling
- Tiling for unstructured-mesh stencils
- Lazy, data-driven compute-communicate
- Runtime code generation
- Multicore graph worklists
- Tensor contractions
- Generalised loop-invariant code motion
- Sparsity in Fourier transforms
- Functional Variational Inference
- Search-based optimisation
- Processor/accelerator microarchitecture, co-design
- MLIR

Technologies

- Finite-volume CFD
- Finite-element
- Finite-difference
- Real-time 3D scene understanding
- Adaptive-mesh CFD
- Contour trees, Reeb graphs
- Unsteady CFD - higher-order flux-reconstruction
- Ab-initio computational chemistry (ONETEP)
- Gaussian belief propagation
- Uncertainty in DNNs
- Near-camera processing
- Quantum computing

Contexts

Automating domain-specific performance optimisations

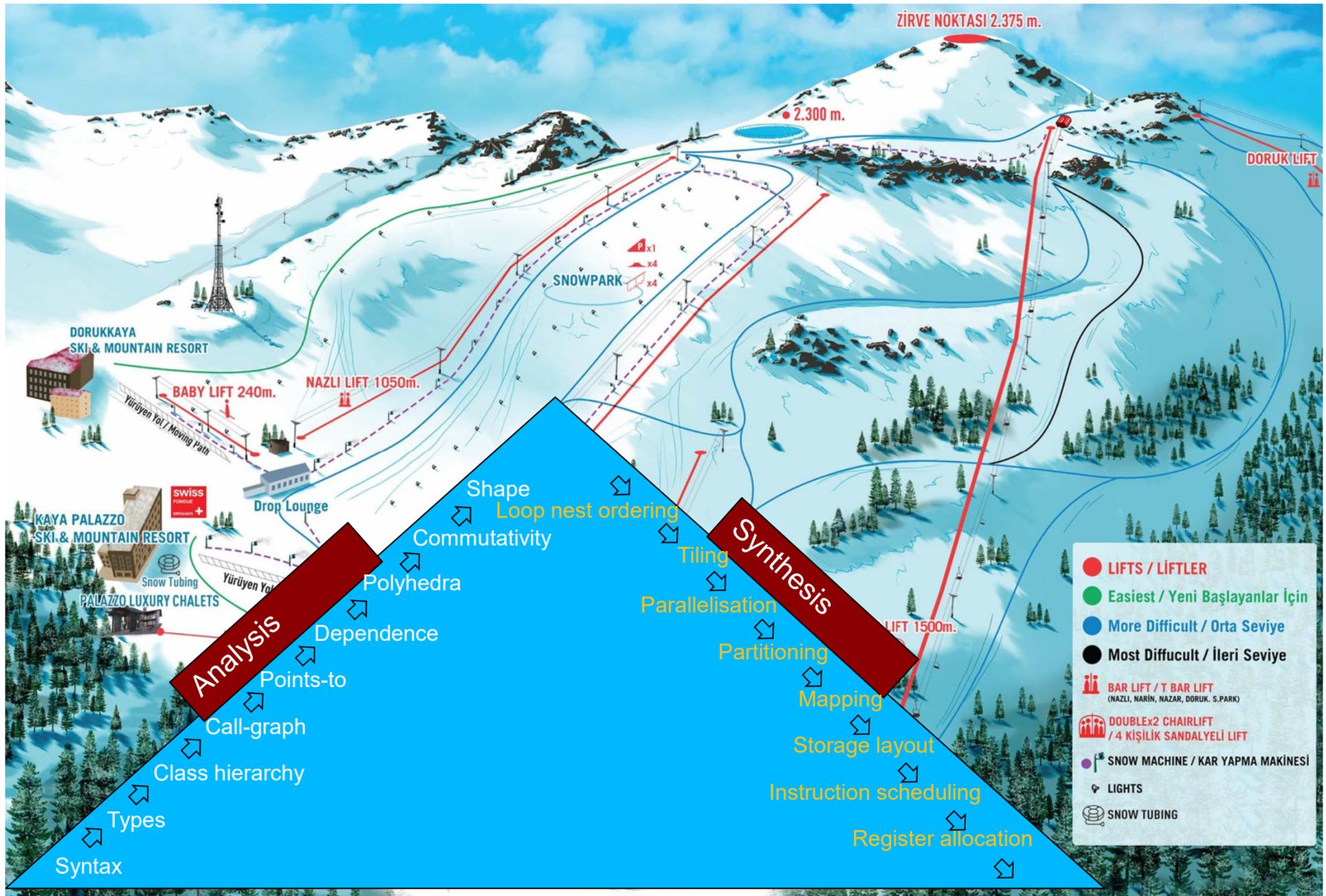
Exploiting higher-level language to get better performance than low level code

- PyOP2/OP2**
Unstructured-mesh stencils
- Firedrake**
Finite-element
- Devito**: finite difference
- SLAMBench**:
3D vision, dense SLAM
- SuperEight**: octree SLAM
- PRAGMaTic**: Unstructured mesh adaptation
- GiMMiK**: small matrix multiply
- TINTL**: Fourier interpolation
- Hypermapper**: design optimisation
- RobotWeb**: distributed localisation

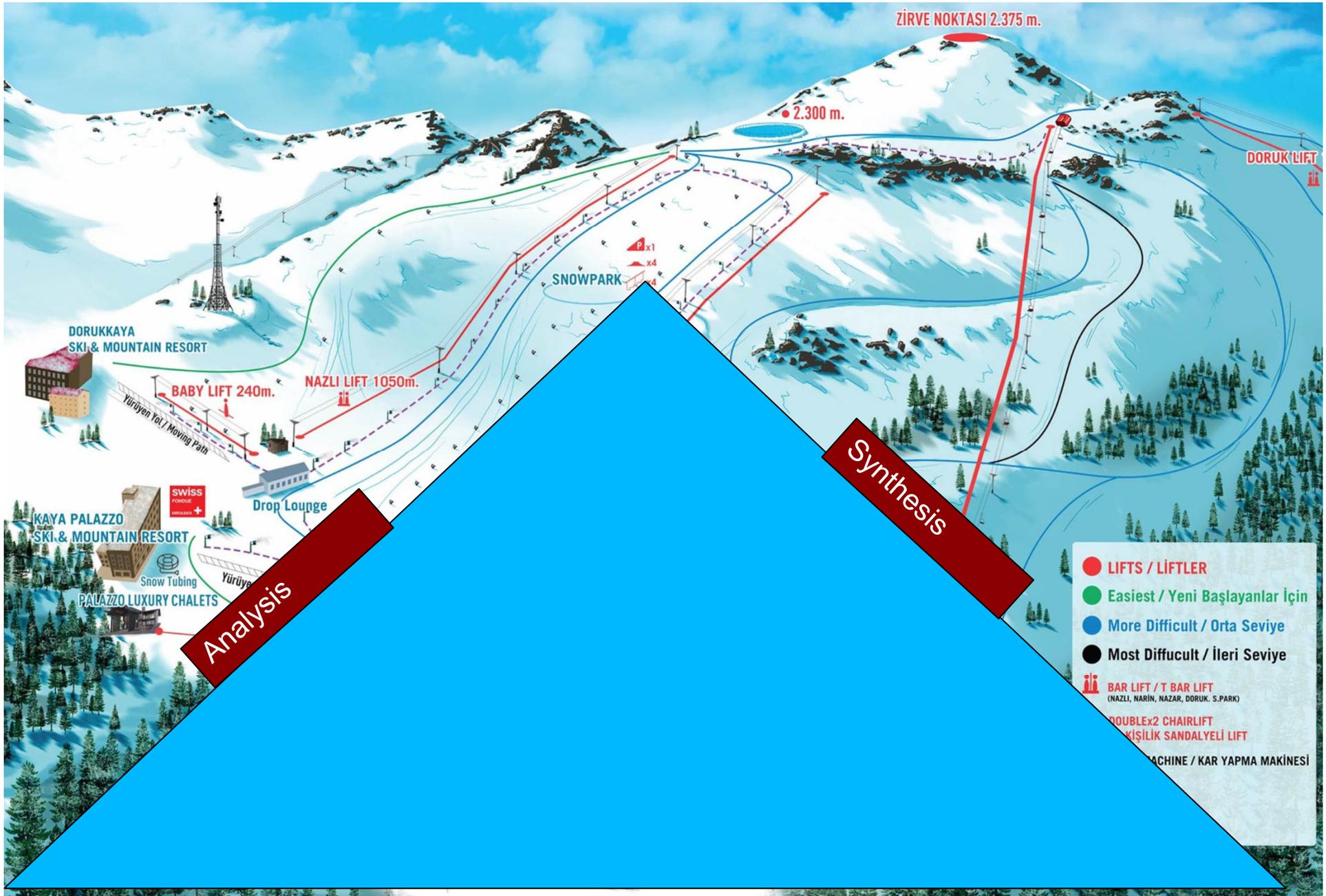
Projects

- Aeroengine turbo-machinery
- Weather and climate
- Glaciers
- Domestic robotics, augmented reality
- Tidal turbine placement
- Formula-1, UAVs, buildings
- Solar energy, drug design
- Medical imaging

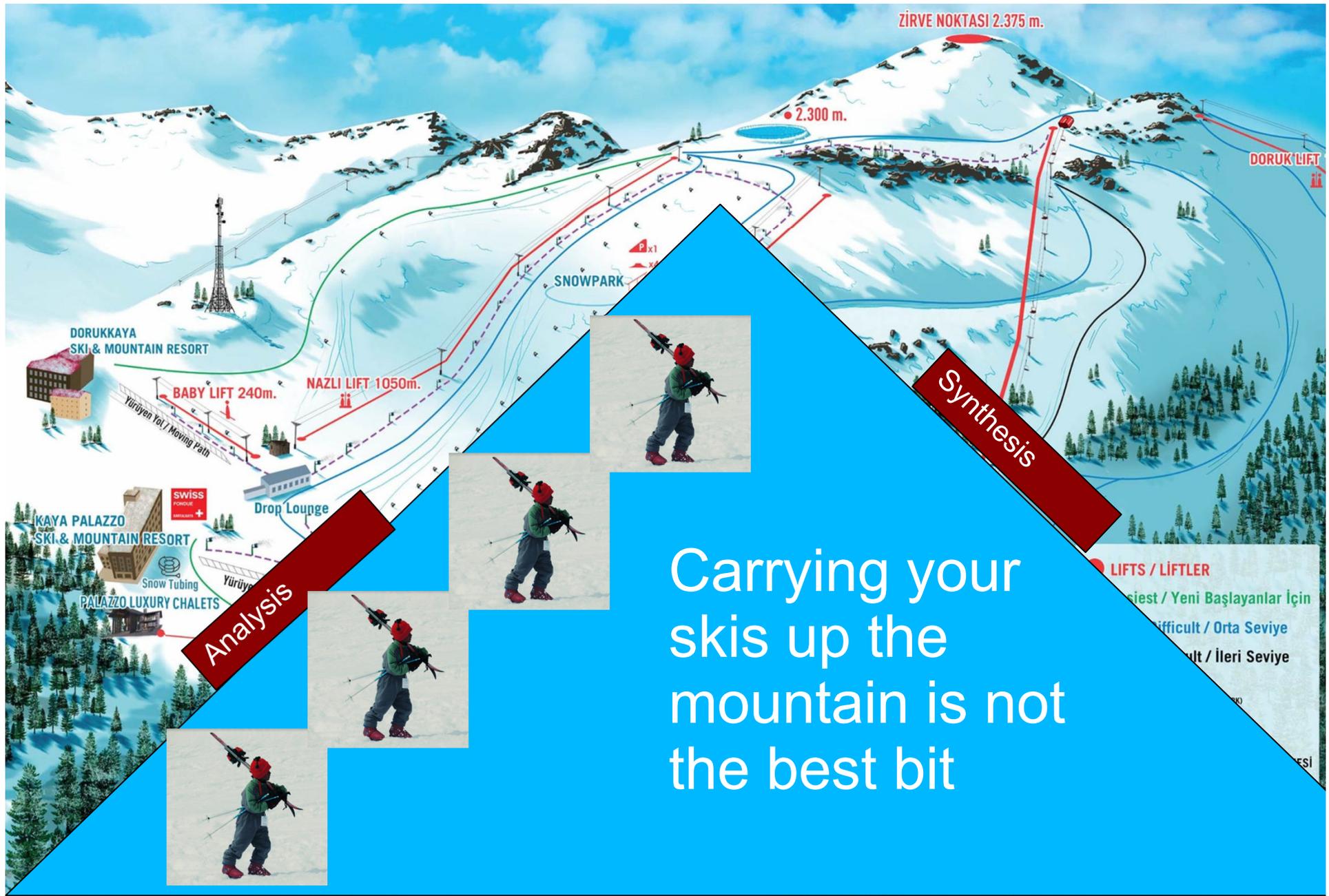
Applications



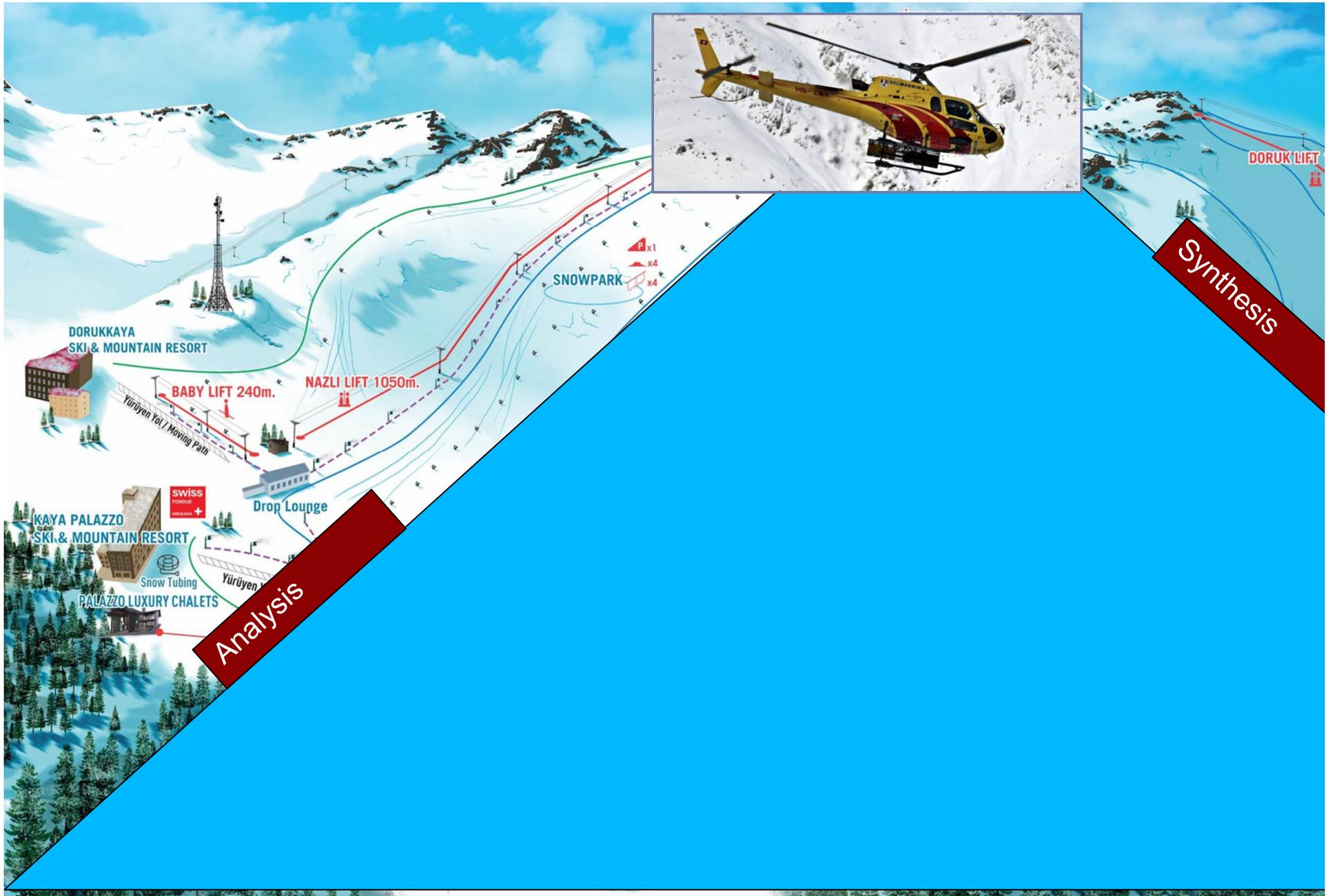
Compilation is like skiing



Compilation is like skiing



Compilation is like skiing



Compilation is like skiing

- Compiling is like skiing
 - Analysis is an uphill struggle
- “Turing Tax”
 - The price you pay for running on a general-purpose computer rather than a specialised one
- What do we call...
 - ***The price you pay for using a general-purpose programming language rather than a DSL?***
- This talk:
 - DSLs really can deliver – my examples: Firedrake, Devito
 - DSL compiler architecture: how do DSLs win?
 - Making the DSL ecosystem work



All of this is Turing Tax!

“Turing tax”: the price we pay for using a general-purpose tool instead of a special-purpose one

Compilation is like skiing

Example DSL:



Firedrake

[Documentation](#) [Download](#) [Team](#) [Citing](#) [Publications](#) [Events](#) [Funding](#) [Contact](#) [GitHub](#) [Jenkins](#)

Firedrake is an automated system for the solution of partial differential equations using the finite element method (FEM). Firedrake uses sophisticated code generation to provide mathematicians, scientists, and engineers with a very high productivity way to create sophisticated high performance simulations.

Features:

- Expressive specification of any PDE using the Unified Form Language from [the FEniCS Project](#).
- Sophisticated, programmable solvers through seamless coupling with [PETSc](#).
- Triangular, quadrilateral, and tetrahedral unstructured meshes.
- Layered meshes of triangular wedges or hexahedra.
- Vast range of finite element spaces.
- Sophisticated automatic optimisation, including sum factorisation for high order elements, and vectorisation.
- Geometric multigrid.
- Customisable operator preconditioners.
- Support for static condensation, hybridisation, and HDG methods.

Latest commits to the Firedrake master branch on Github

Merge pull request #1520 from firedrakeproject/wence/feature/assemble-diagonal

Lawrence Mitchell authored at 22/10/2019, 09:14:34

tests: Check that getting diagonal of matrix works

Lawrence Mitchell authored at 21/10/2019, 13:04:04

matfree: Add getDiagonal method to implicit matrices

Lawrence Mitchell authored at 18/10/2019, 10:19:48

assemble: Add option to assemble diagonal of 2-form into Dat

Lawrence Mitchell authored at 18/10/2019, 10:08:37

Merge pull request #1509 from firedrakeproject/wence/patch-c-wrapper

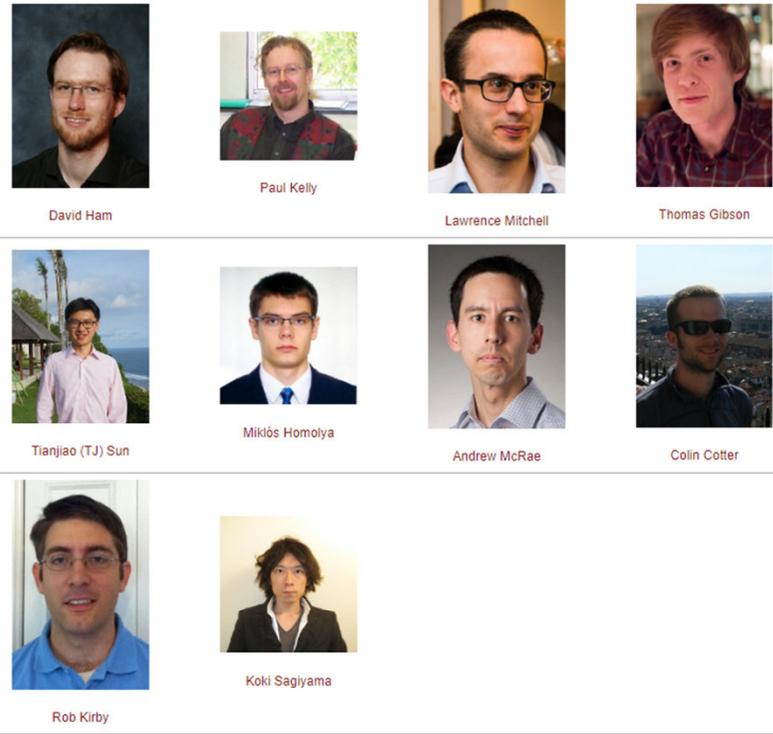
The screenshot shows the website for the Firedrake project. At the top left is a red dragon logo. To its right is the word "Firedrake" in a large, stylized font. Below the logo and name is a navigation bar with links for "Documentation", "Download", "Team", and "Citir". A red arrow points to the "Team" link. The main content area on the left contains a paragraph describing Firedrake as an automated system for solving partial differential equations using the finite element method (FEM). Below this is a "Features:" section with a bulleted list of capabilities. On the right side of the page, there are two sections: "Active team members" and "Former team members", each displaying a grid of portrait photos of team members with their names listed below.

Firedrake is an automated system for the solution of partial differential equations using the finite element method (FEM). Firedrake uses sophisticated code generation to help mathematicians, scientists, and engineers with a very high productivity way to perform sophisticated high performance simulations.

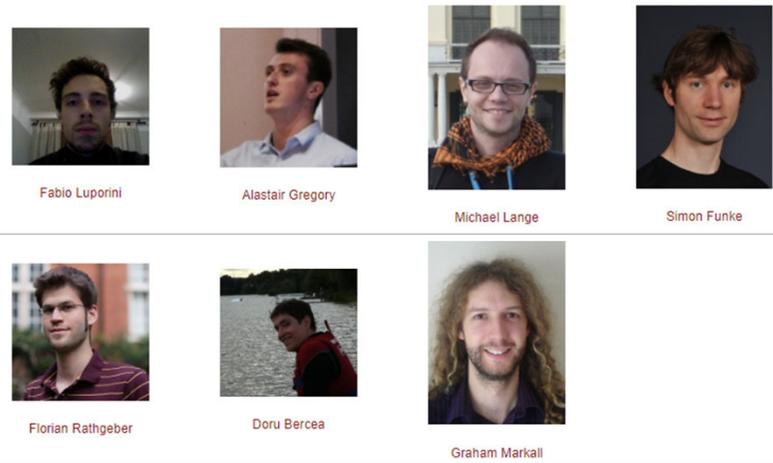
Features:

- Expressive specification of any PDE using the Unified Form Language [Project](#).
- Sophisticated, programmable solvers through seamless coupling with
- Triangular, quadrilateral, and tetrahedral unstructured meshes.
- Layered meshes of triangular wedges or hexahedra.
- Vast range of finite element spaces.
- Sophisticated automatic optimisation, including sum factorisation for elements, and vectorisation.
- Geometric multigrid.
- Customisable operator preconditioners.
- Support for static condensation, hybridisation, and HDG methods.

Active team members



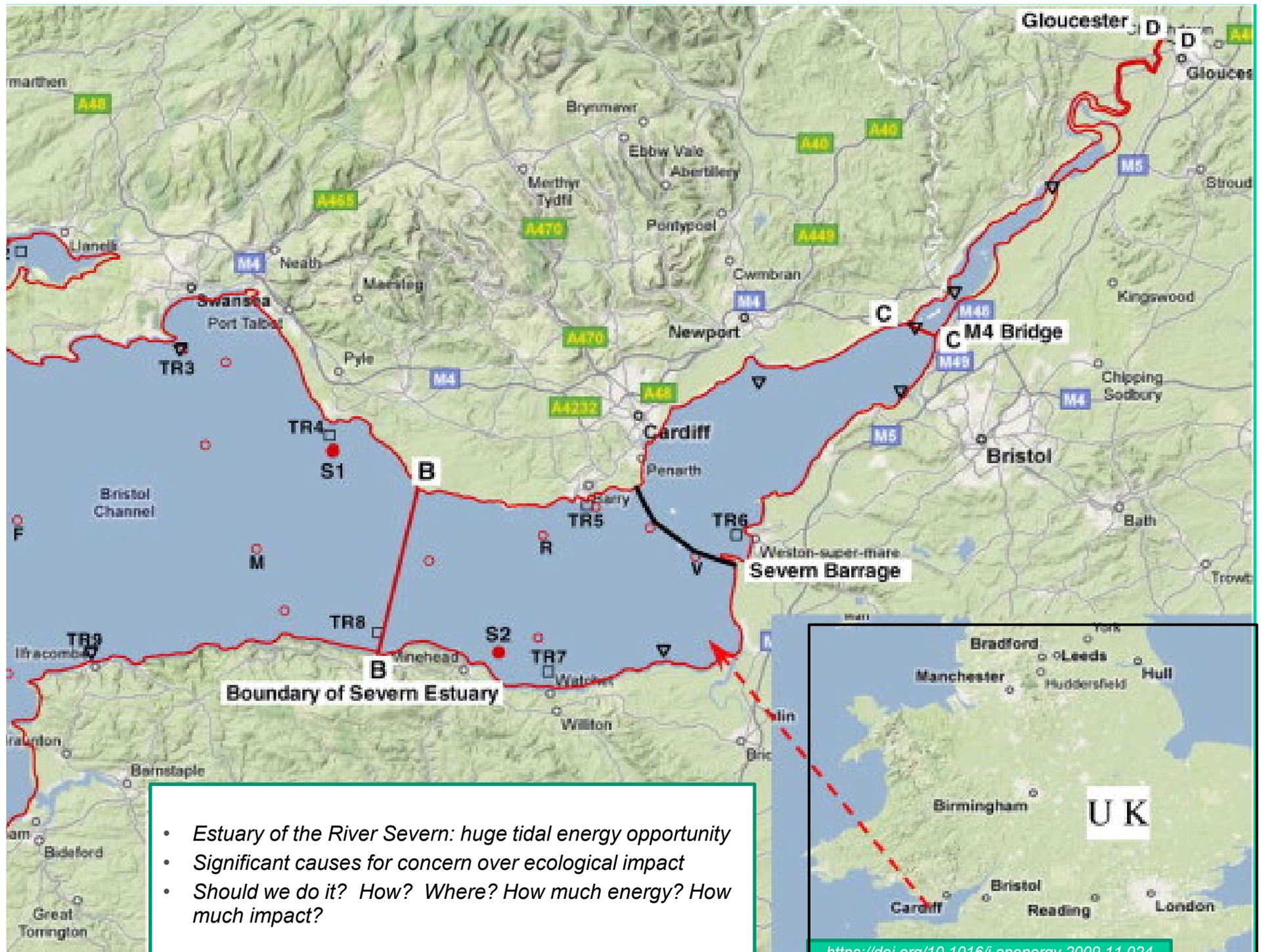
Former team members



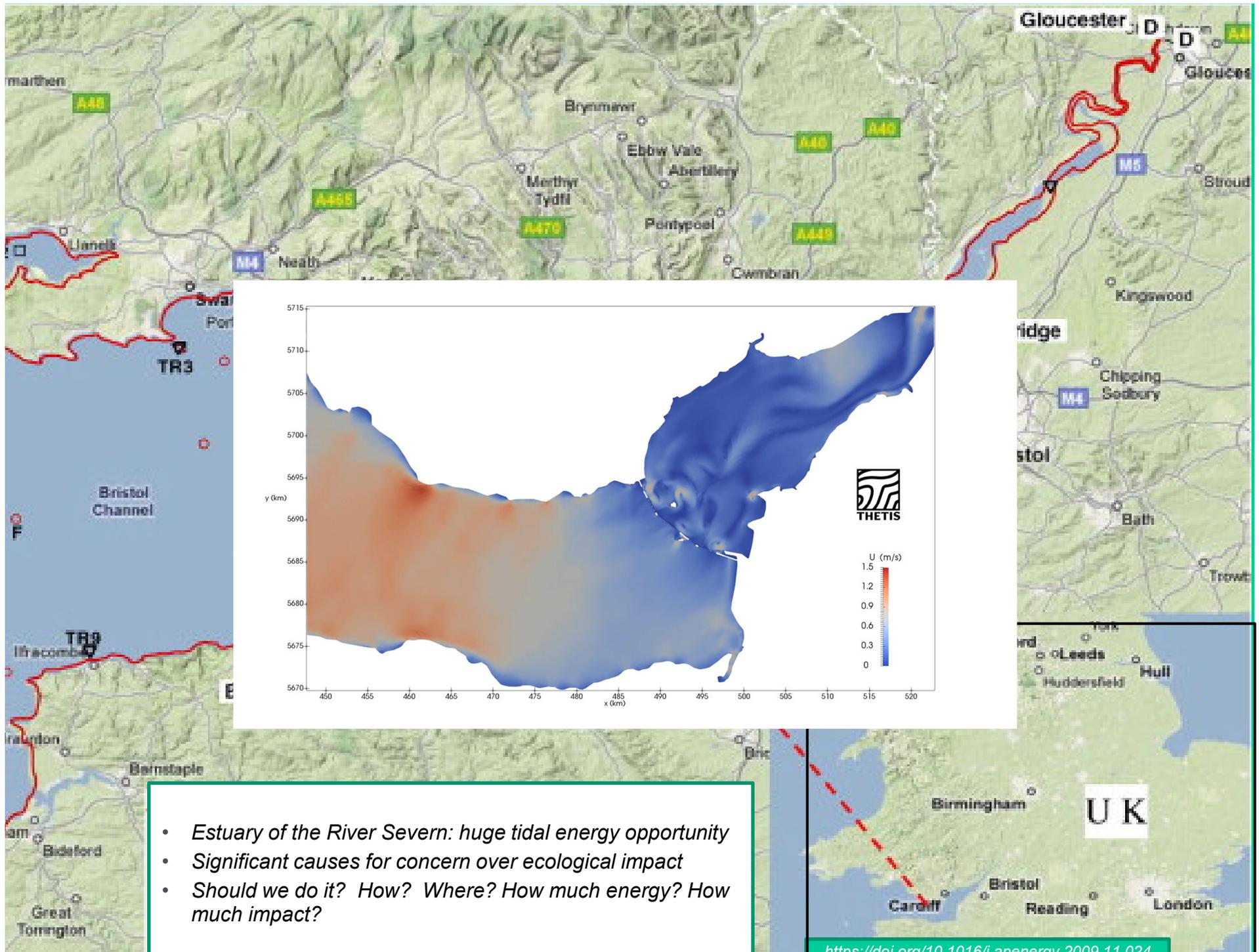
- Firedrake is used in:
 - **Thetis:** unstructured grid coastal modelling framework

The screenshot shows the website for the Thetis project at thetisproject.org. The page features a navigation bar with links for Documentation, Download, Team, Publications, Funding, Contact, GitHub, and Jenkins. The main content area is titled "The Thetis project" and describes it as an unstructured grid coastal ocean model built using the Firedrake finite element framework. It mentions that Thetis currently consists of 2D depth-averaged and full 3D baroclinic models. A "Current development status" box indicates the latest status is "build passing" and notes that the source code is hosted on GitHub and tested using Jenkins. Below the text are four video thumbnails: "Idealized river plume simulation", "Baroclinic eddies test case", "Thetis Tidal Barrage simulation", and "Thetis Two Lagoon Simulation".

28 ■ What is it used for? By whom?

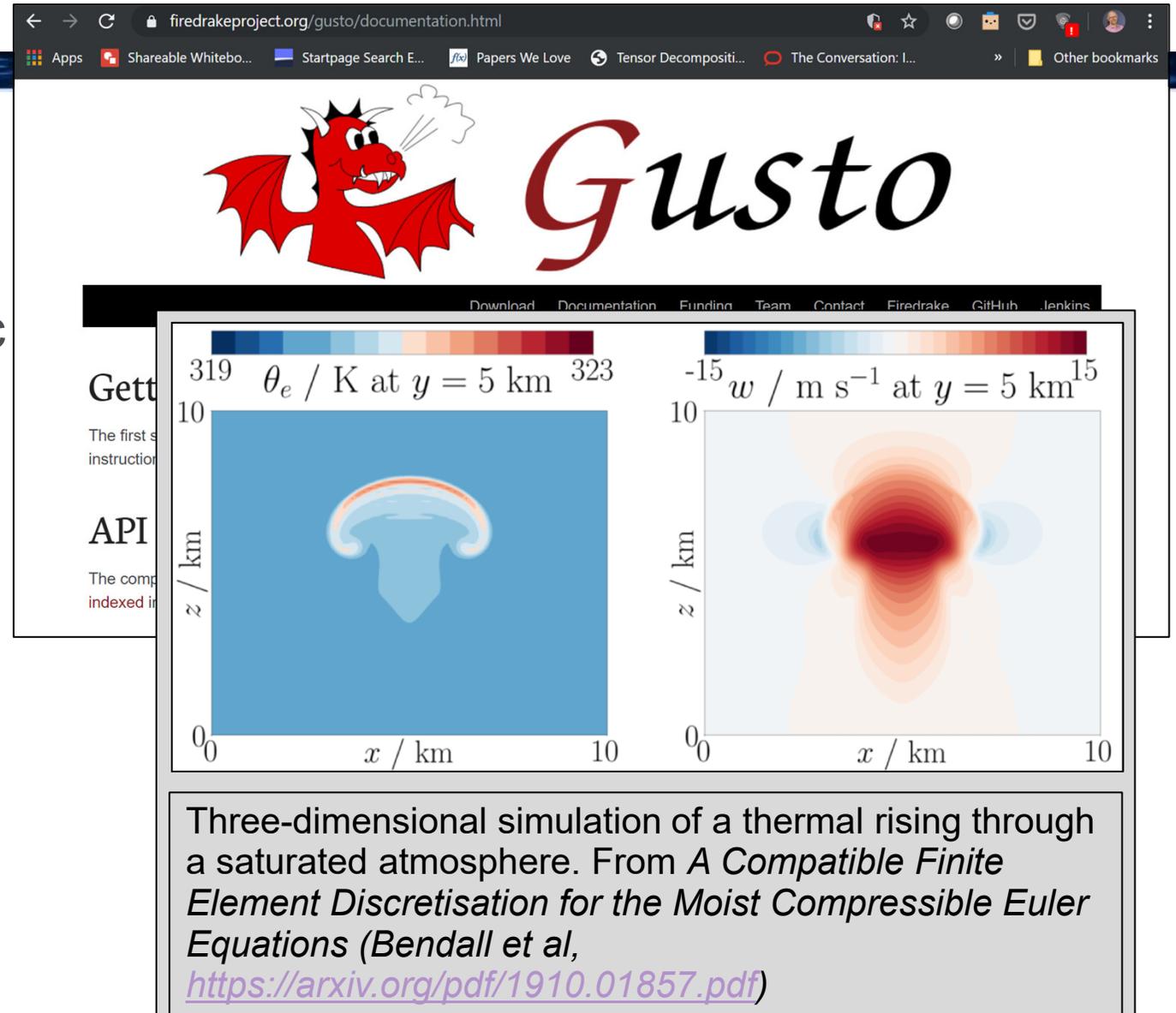


- *Estuary of the River Severn: huge tidal energy opportunity*
- *Significant causes for concern over ecological impact*
- *Should we do it? How? Where? How much energy? How much impact?*



- *Estuary of the River Sever: huge tidal energy opportunity*
- *Significant causes for concern over ecological impact*
- *Should we do it? How? Where? How much energy? How much impact?*

- Firedrake is used in:
 - **Gusto:** atmospheric modelling framework being used to prototype the next generation of weather and climate simulations for the UK Met Office



The screenshot shows the website `firedrakeproject.org/gusto/documentation.html`. At the top, there is a red dragon logo and the word "Gusto" in a large, stylized font. Below this, there are navigation links: "Download", "Documentation", "Funding", "Team", "Contact", "Firedrake", "GitHub", and "Jenkins". The main content area features two 2D contour plots side-by-side, both with axes labeled x / km and z / km ranging from 0 to 10. The left plot is titled θ_e / K at $y = 5 \text{ km}$ and shows a color scale from 319 to 323. The right plot is titled $w / \text{m s}^{-1}$ at $y = 5 \text{ km}$ and shows a color scale from -15 to 15. Both plots show a central region of high values (red/orange) that tapers off as it moves away from the center, representing a thermal rising through a saturated atmosphere. Below the plots, there is a text box containing the following information:

Three-dimensional simulation of a thermal rising through a saturated atmosphere. From *A Compatible Finite Element Discretisation for the Moist Compressible Euler Equations* (Bendall et al, <https://arxiv.org/pdf/1910.01857.pdf>)

■ Firedrake is used in:

■ **Icepack**: a framework for modeling the flow of glaciers and ice sheets, developed at the Polar Science Center at the University of Washington

icepack 0.0.3

Search docs

BASICS

- Overview
- Background
- Installation
- Contact

TUTORIALS

- Meshes, functions
- Synthetic ice shelf
- Larsen Ice Shelf
- Synthetic ice stream
- Inverse problems
- Ice streams, once more

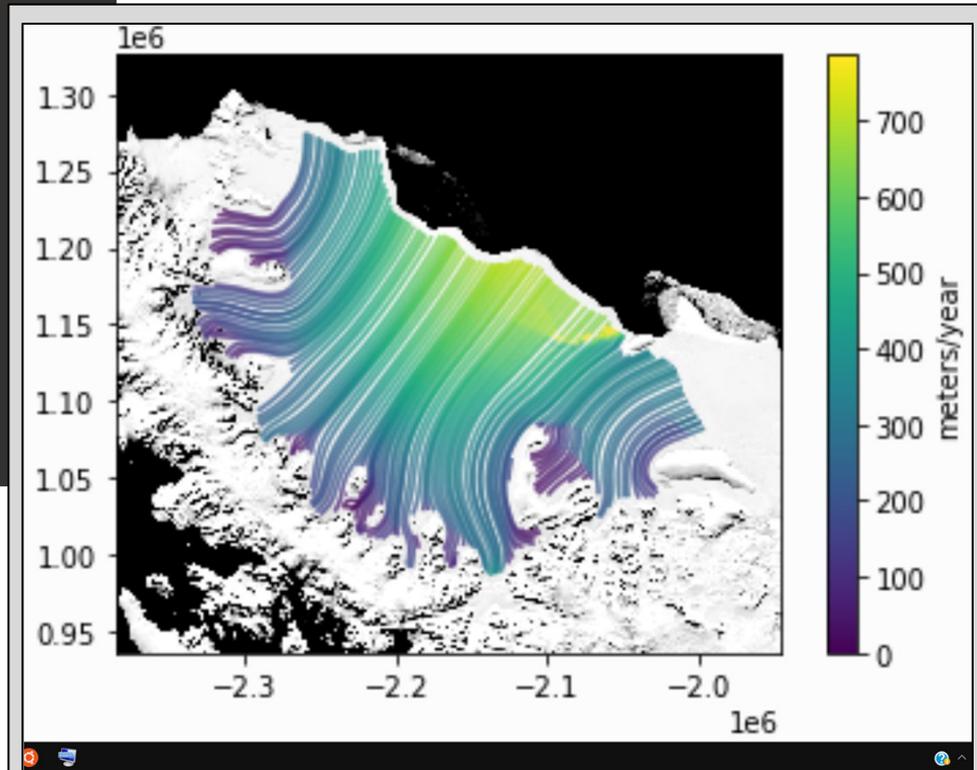
DEVELOPMENT

- Contributing
- Testing

Docs » icepack [View page source](#)

icepack

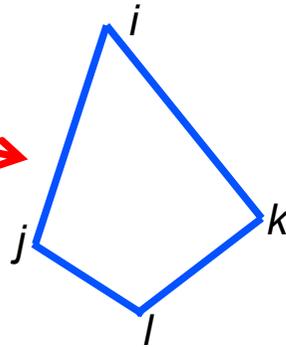
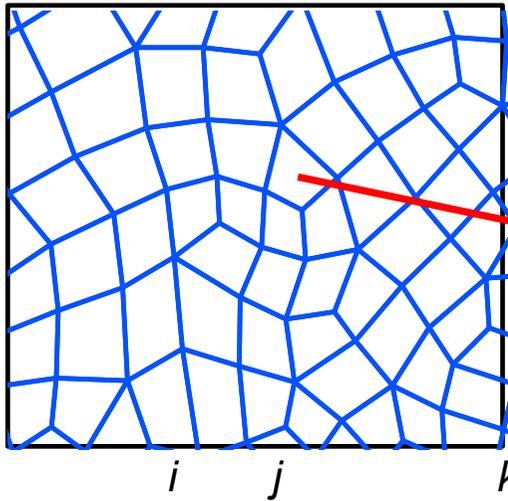
Welcome to the documentation for *icepack*, a python library for modeling the flow of ice sheets and glaciers! The main design goals for icepack are:



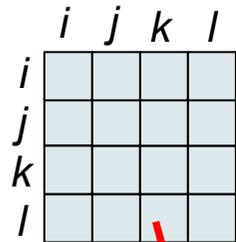
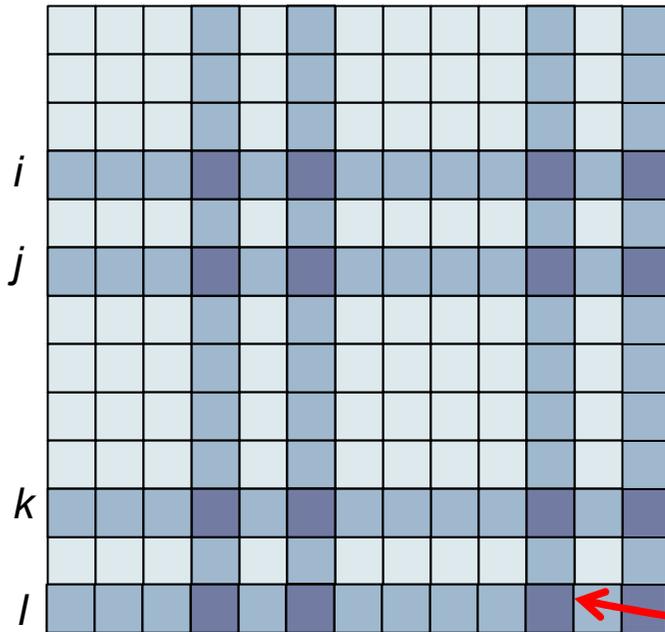
Larsen ice shelf model, from the Icepack tutorial by Daniel Shapero (<https://icepack.github.io/icepack.demo.02-larsen-ice-shelf.html>)

■ What is it used for? By whom?

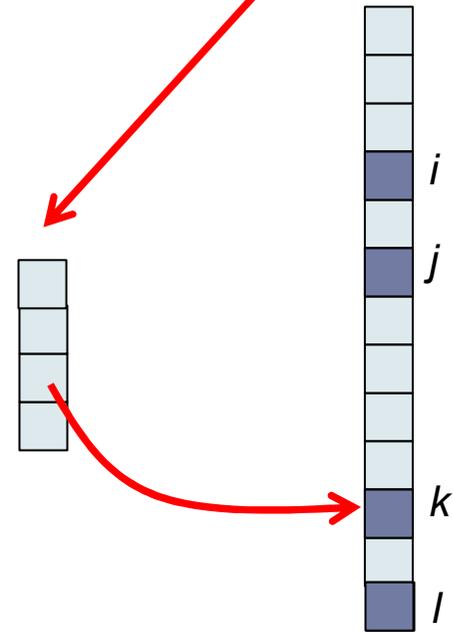
The finite element method in outline



```
do element = 1, N
  assemble(element):
     $\int_{\Omega} vL(u^{\delta})dX = \int_{\Omega} vqdX.$ 
end do
```



$$Ax = b$$



- Key data structures: Mesh, dense local assembly matrices, sparse global system matrix, and RHS vector

■ Local assembly:

- Computes local assembly matrix
- Using:
 - The (weak form of the) PDE
 - The discretisation
- Key operation is evaluation of expressions over basis function representation of the element

■ Mesh traversal:

- *PyOP2*
- *Loops over the mesh*
- *Key is orchestration of data movement*

■ Solver:

- Interfaces to standard solvers through PetSc

Example: Burgers equation

$$\int_{\Omega} \frac{u^{n+1} - u^n}{dt} \cdot v + ((u^{n+1} \cdot \nabla)u^{n+1}) \cdot v + \nu \nabla u^{n+1} \cdot \nabla v \, dx = 0.$$

- From the weak form of the PDE, we derive an equation to solve, that determines the state at each timestep in terms of the previous timestep

- Transcribe into Python – u is u^{n+1} , u_* is u^n :

```
F = (inner((u - u_)/timestep, v)
      + inner(dot(u,nabla_grad(u)), v) + nu*inner(grad(u), grad(v)))*dx
```

- Set up the equation and solve for the next timestep u :
`solve(F == 0, u)`
- At this point, Firedrake generates code to assemble a linear system, runs it and calls a linear solver (we use PETSC)

Burgers equation

```
from firedrake import *
n = 50
mesh = UnitSquareMesh(n, n)

# We choose degree 2 continuous Lagrange polynomials. We also need a
# piecewise linear space for output purposes::

V = VectorFunctionSpace(mesh, "CG", 2)
V_out = VectorFunctionSpace(mesh, "CG", 1)

# We also need solution functions for the current and the next timestep::

u_ = Function(V, name="Velocity")
u = Function(V, name="VelocityNext")

v = TestFunction(V)

# We supply an initial condition::

x = SpatialCoordinate(mesh)
ic = project(as_vector([sin(pi*x[0]), 0]), V)

# Start with current value of u set to the initial condition, and use the
# initial condition as our starting guess for the next value of u::

u_.assign(ic)
u.assign(ic)

# nu :math:`\nu` is set to a (fairly arbitrary) small constant value::

nu = 0.0001

timestep = 1.0/n

# Define the residual of the equation::

F = (inner((u - u_)/timestep, v)
      + inner(dot(u, nabla_grad(u)), v) + nu*inner(grad(u), grad(v)))*dx

outfile = File("burgers.pvd")

outfile.write(project(u, V_out, name="Velocity"))

# Finally, we loop over the timesteps solving the equation each time::

t = 0.0
end = 0.5
while (t <= end):
    solve(F == 0, u)
    u_.assign(u)
    t += timestep
    outfile.write(project(u, V_out, name="Velocity"))
```

- Firedrake implements the Unified Form Language (UFL)
- Embedded in Python

$$\int_{\Omega} \frac{u^{n+1} - u^n}{dt} \cdot v + ((u^{n+1} \cdot \nabla)u^{n+1}) \cdot v + \nu \nabla u^{n+1} \cdot \nabla v \, dx = 0.$$

- From the weak form of the PDE, we derive an equation to solve, that determines the state at each timestep in terms of the previous timestep
- Transcribe into Python – u is u^{n+1} , $u_$ is u^n :

```
F = (inner((u - u_)/timestep, v)
      + inner(dot(u, nabla_grad(u)), v) + nu*inner(grad(u), grad(v)))*dx
```

- (UFL is also the DSL of the FEniCS project)

- What does its DSL actually look like?

```
from firedrake import *
n = 50
mesh = UnitSquareMesh(n, n)
```

```
mesh = UnitSquareMesh(n, n)
```

```
# We choose degree 2 continuous Lagrange polynomials.
# piecewise linear space for output purposes::
```

```
V = VectorFunctionSpace(mesh, "CG", 2)
V_out = VectorFunctionSpace(mesh, "CG", 1)
```

```
V = VectorFunctionSpace(mesh, "CG", 2)
V_out = VectorFunctionSpace(mesh, "CG", 1)
```

```
# We also need solution functions for the current and the next timestep::
```

```
u_ = Function(V, name="Velocity")
u = Function(V, name="VelocityNext")
```

```
u_ = Function(V, name="Velocity")
u = Function(V, name="VelocityNext")
```

```
v = TestFunction(V)
```

```
# We supply an initial condition::
```

```
# set up initial conditions for u and u_
```

```
x = SpatialCoordinate(mesh)
ic = project(as_vector([sin(pi*x[0]), 0]), V)
```

```
# Start with current value of u set to
# initial condition as our starting guess
```

```
# Define the residual of the equation::
```

```
u_.assign(ic)
u.assign(ic)
```

```
F = (inner((u - u_)/timestep, v)
      + inner(dot(u, nabla_grad(u)), v) + nu*inner(grad(u), grad(v)))*dx
```

```
# :math:`\nu` is set to a (fairly arbitrary)
```

```
nu = 0.000
```

```
timestep =
```

```
# Define t
```

```
F = (inner
      + inn
```

```
outfile =
```

```
outfile.wr
```

```
# Finally,
```

```
t = 0.0
```

```
end = 0.5
```

```
while (t <
```

```
    solve(F == 0, u)
```

```
    u_.assign(u)
```

```
    t += timestep
```

```
    outfile.write(project(u, V_out, name="Velocity"))
```

```
t = 0.0
```

```
end = 0.5
```

```
while (t <= end):
```

```
    solve(F == 0, u)
```

```
    u_.assign(u)
```

```
    t += timestep
```

```
    outfile.write(project(u, V_out, name="Velocity"))
```

■ What does its DSL actually look like?

```

#include <math.h>
#include <petsc.h>

void wrap_form00_cell_integral_otherwise(int const start, int const end, Mat const mat0, double const *__restrict__ dat1, double const *__restrict__ dat0, int const *__restrict__ map0, int const *__restrict__ map1)
{
  double form_t0...t16;
  double const form_t17[7] = { ... };
  double const form_t18[7 * 6] = { ... };
  double const form_t19[7 * 6] = { ... };
  double form_t2;
  double const form_t20[7 * 6] = { ... };
  double form_t21...t37;
  double form_t38[6];
  double form_t39[6];
  double form_t4;
  double form_t40...t45;
  double form_t5...t9;
  double t0[6 * 2];
  double t1[3 * 2];
  double t2[6 * 2 * 6 * 2];

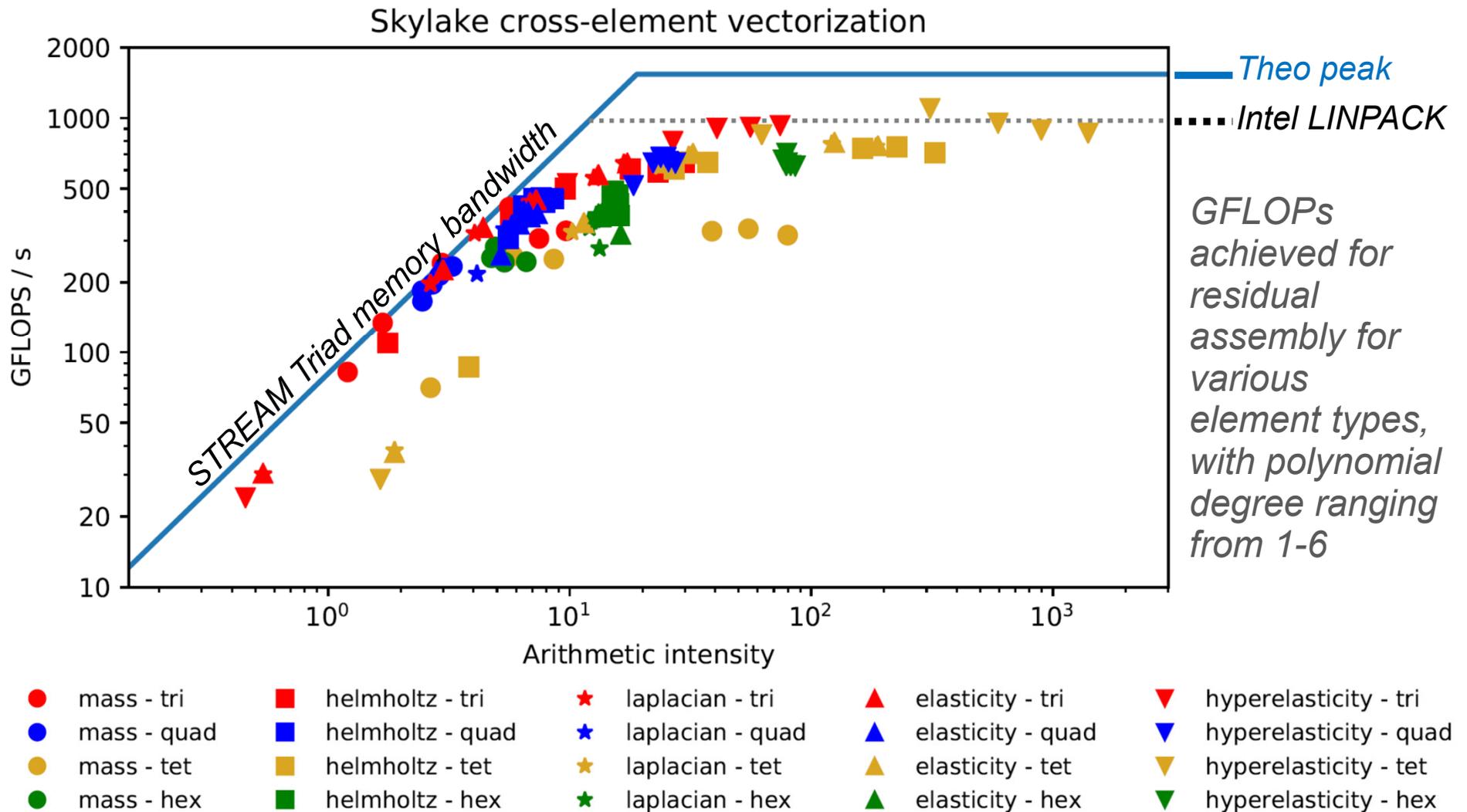
  for (int n = start; n <= -1 + end; ++n)
  {
    for (int i4 = 0; i4 <= 5; ++i4)
      for (int i5 = 0; i5 <= 1; ++i5)
        for (int i6 = 0; i6 <= 5; ++i6)
          for (int i7 = 0; i7 <= 1; ++i7)
            t2[24 * i4 + i2 * 15 + 2 * i6 + i7] = 0.0;
    for (int i2 = 0; i2 <= 2; ++i2)
      for (int i3 = 0; i3 <= 1; ++i3)
        t1[2 * i2 + i3] = dat1[2 * map1[3 * n + i2] + i3];
    for (int i0 = 0; i0 <= 5; ++i0)
      for (int i1 = 0; i1 <= 1; ++i1)
        t0[2 * i0 + i1] = dat0[2 * map0[6 * n + i0] + i1];
    form_t0 = -1.0 * t1[1];
    form_t1 = form_t0 + t1[3];
    form_t2 = -1.0 * t1[0];
    form_t3 = form_t2 + t1[2];
    form_t4 = form_t0 + t1[5];
    form_t5 = form_t2 + t1[4];
    form_t6 = form_t3 * form_t4 + -1.0 * form_t5 * form_t1;
    form_t7 = 1.0 / form_t6;
    form_t8 = form_t7 * -1.0 * form_t1;
    form_t9 = form_t4 * form_t7;
    form_t10 = form_t3 * form_t7;
    form_t11 = form_t7 * -1.0 * form_t5;
    form_t12 = 0.0001 * (form_t8 * form_t9 + form_t10 * form_t11);
    form_t13 = 0.0001 * (form_t8 * form_t8 + form_t10 * form_t10);
    form_t14 = 0.0001 * (form_t9 * form_t9 + form_t11 * form_t11);
    form_t15 = 0.0001 * (form_t9 * form_t8 + form_t11 * form_t10);
    form_t16 = fabs(form_t6);
    for (int form_ip = 0; form_ip <= 6; ++form_ip)
    {
      form_t26 = 0.0; form_t25 = 0.0; form_t24 = 0.0; form_t23 = 0.0; form_t22 = 0.0; form_t21 = 0.0;
      for (int form_i = 0; form_i <= 5; ++form_i)
      {
        form_t21 = form_t21 + form_t20[6 * form_ip + form_i] * t0[1 + 2 * form_i];
        form_t22 = form_t22 + form_t19[6 * form_ip + form_i] * t0[1 + 2 * form_i];
        form_t23 = form_t23 + form_t20[6 * form_ip + form_i] * t0[2 * form_i];
        form_t24 = form_t24 + form_t19[6 * form_ip + form_i] * t0[2 * form_i];
        form_t25 = form_t25 + form_t18[6 * form_ip + form_i] * t0[1 + 2 * form_i];
        form_t26 = form_t26 + form_t18[6 * form_ip + form_i] * t0[2 * form_i];
      }
      form_t27 = form_t17[form_ip] * form_t16;
      form_t28 = form_t27 * form_t15;
      form_t29 = form_t27 * form_t14;
      form_t30 = form_t27 * (form_t26 * form_t9 + form_t25 * form_t11);
      form_t31 = form_t27 * form_t13;
      form_t32 = form_t27 * form_t12;
      form_t33 = form_t27 * (form_t26 * form_t8 + form_t25 * form_t10);
      form_t34 = form_t27 * (form_t11 * form_t24 + form_t10 * form_t23);
      form_t35 = form_t27 * (form_t9 * form_t22 + form_t8 * form_t21);
      form_t36 = form_t27 * (50.0 + form_t9 * form_t24 + form_t8 * form_t23);
      form_t37 = form_t27 * (50.0 + form_t11 * form_t22 + form_t10 * form_t21);
      for (int form_k0 = 0; form_k0 <= 5; ++form_k0)
      {
        form_t38[form_k0] = form_t18[6 * form_ip + form_k0] * form_t37;
        form_t39[form_k0] = form_t18[6 * form_ip + form_k0] * form_t36;
      }
      for (int form_j0 = 0; form_j0 <= 5; ++form_j0)
      {
        form_t40 = form_t18[6 * form_ip + form_j0] * form_t35;
        form_t41 = form_t18[6 * form_ip + form_j0] * form_t34;
        form_t42 = form_t20[6 * form_ip + form_j0] * form_t31 + form_t18[6 * form_ip + form_j0] * form_t33 + form_t19[6 * form_ip + form_j0] * form_t32;
        form_t43 = form_t20[6 * form_ip + form_j0] * form_t28 + form_t18[6 * form_ip + form_j0] * form_t30 + form_t19[6 * form_ip + form_j0] * form_t29;
        for (int form_k0_0 = 0; form_k0_0 <= 5; ++form_k0_0)
        {
          form_t44 = form_t43 * form_t19[6 * form_ip + form_k0_0];
          form_t45 = form_t42 * form_t20[6 * form_ip + form_k0_0];
          t2[24 * form_j0 + 2 * form_k0_0] = t2[24 * form_j0 + 2 * form_k0_0] + form_t45 + form_t18[6 * form_ip + form_j0] * form_t39[form_k0_0] + form_t44;
          t2[13 + 24 * form_j0 + 2 * form_k0_0] = t2[13 + 24 * form_j0 + 2 * form_k0_0] + form_t45 * form_t18[6 * form_ip + form_j0] * form_t41;
          t2[1 + 24 * form_j0 + 2 * form_k0_0] = t2[1 + 24 * form_j0 + 2 * form_k0_0] + form_t18[6 * form_ip + form_k0_0] * form_t41;
          t2[12 + 24 * form_j0 + 2 * form_k0_0] = t2[12 + 24 * form_j0 + 2 * form_k0_0] + form_t18[6 * form_ip + form_k0_0] * form_t40;
        }
      }
    }
    MatSetValuesBlockedLocal(mat0, 6, &(map0[6 * n]), 6, &(map0[6 * n]), &(t2[0]), ADD_VALUES);
  }
}

```

- Generated code to assemble the resulting linear system matrix
- Executed at each triangle in the mesh
- Accesses degrees of freedom shared with neighbour triangles through indirection map

Firedrake: single-node AVX512 performance

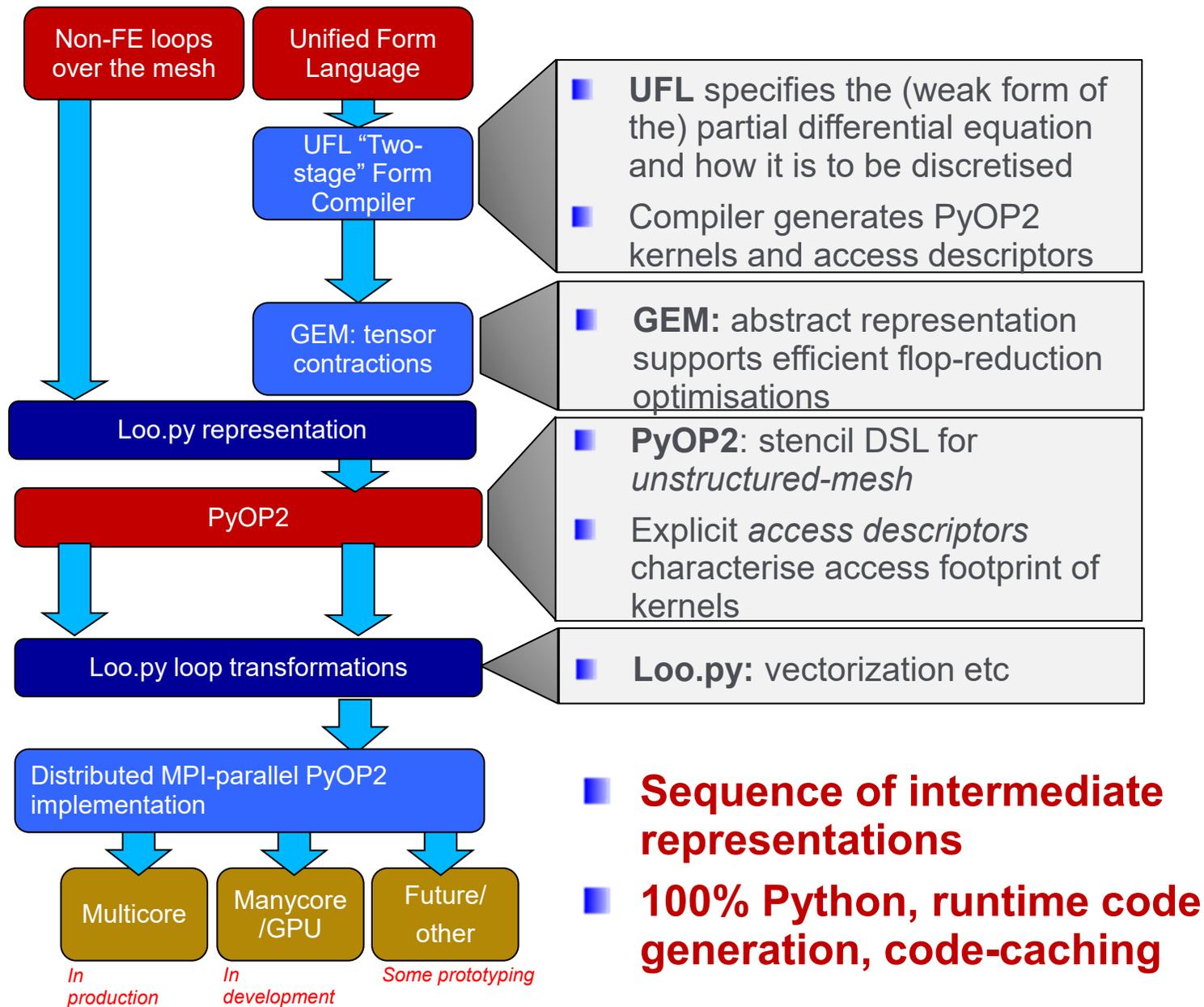
■ Does it generate good code?



[Skylake Xeon Gold 6130 (on all 16 cores, 2.1GHz, turboboost off, Stream: 36.6GB/s, GCC7.3 -march=native)]

A study of vectorization for matrix-free finite element methods, Tianjiao Sun et al
 IJHPCA 2020 <https://arxiv.org/abs/1903.08243>

Firedrake: compiler architecture



Another example DSL:

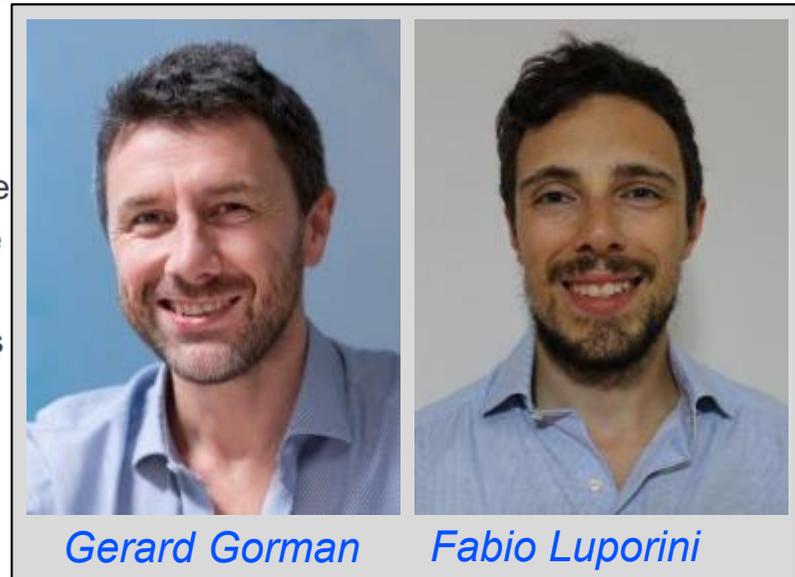
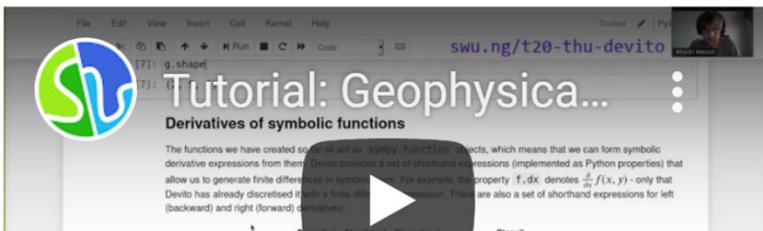


Devito: Symbolic Finite Difference Computation

Devito is a domain-specific Language (DSL) and code generation framework for the design of highly optimised finite difference kernels for use in inversion methods. Devito utilises [SymPy](#) to allow the definition of operators from high-level symbolic equations and generates optimised and automatically tuned code specific to a given target architecture.

Symbolic computation is a powerful tool that allows users to:

- Build complex solvers from only a few lines of high-level code
- Use automated performance optimisation for generated code
- Adjust stencil discretisation at runtime as required
- (Re-)development of solver code in hours rather than months



Gerard Gorman

Fabio Luporini

And many many more!

- Devito automates the finite difference method for solving PDEs
- Widely used for fluid dynamics, wave propagation
- Devito is mostly used to solve inversion problems
 - Use automatic differentiation of the solver
 - To solve for the conditions that explain the observations
 - “Full Waveform Inversion” (FWI)
- **Seismic inversion**
 - Understand geological structures from reflected sound waves
- **Ultrasound imaging of the brain**
 - Diagnose brain injuries from ultrasound transmission

Define the wavefield from model setup.

```
u = TimeFunc(time_order=2, space_order=2)
```

Write down the acoustic wave PDE:

```
pde = model.m*u.dt2 - u.laplace + model.damp*u.dt
```

Solve by time-marching:

```
stencil = Eq(u.forward, solve(pde, u.forward))
```

Define source injection and receiver:

```
src_term = src.inject(field=u.forward, pr=src*dt**2/mo
```

```
rec_term = rec.interpolate(expr=u.forward)
```

Generate code for the timestepping operator:

```
op = Operator([stencil] + src_term + rec_term,
              subs=model.spacing_map)
```

Run code (MPI+GPU), to yield receiver values:

```
op(time=time_range.num-1, dt=model.critical_dt)
```

Acoustic wave equation, with damping:

$$\begin{cases} m \frac{d^2 u(x,t)}{dt^2} - \nabla^2 u(x,t) + \eta \frac{du(x,t)}{dt} = q \text{ in } \Omega \\ u(\cdot, 0) = 0 \\ \frac{du(x,t)}{dt} \Big|_{t=0} = 0 \end{cases}$$

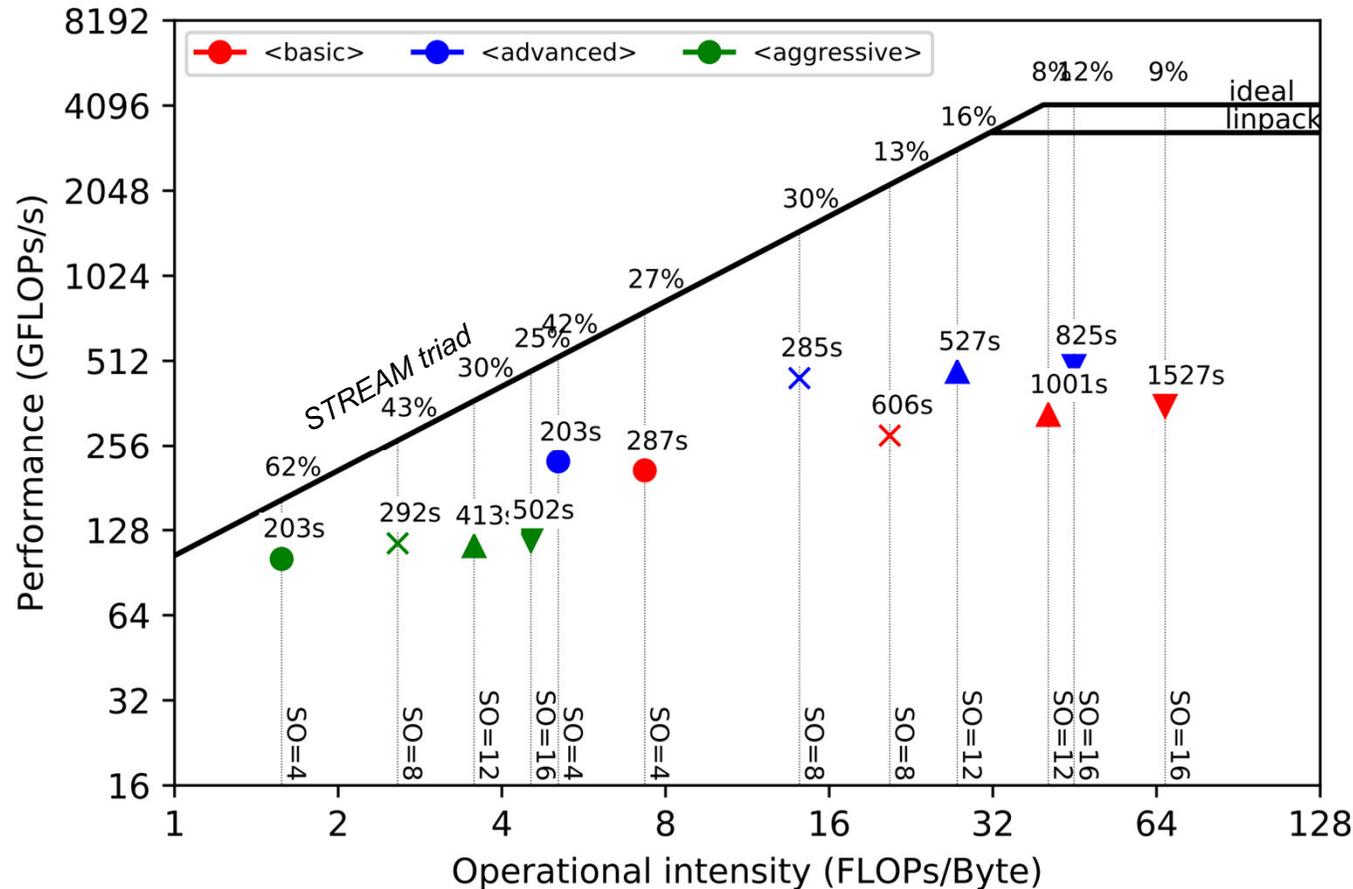
We inject initial sound wave at source point, and monitor the signal at a receiver.

We derive and generate the stencil operator code, then run it a specified number of timesteps

Slightly simplified from:

https://slimgroup.github.io/Devito-Examples/tutorials/01_modelling/

Code at this basic level of abstraction is in production, at scale, running at multiple petaflops 24/7



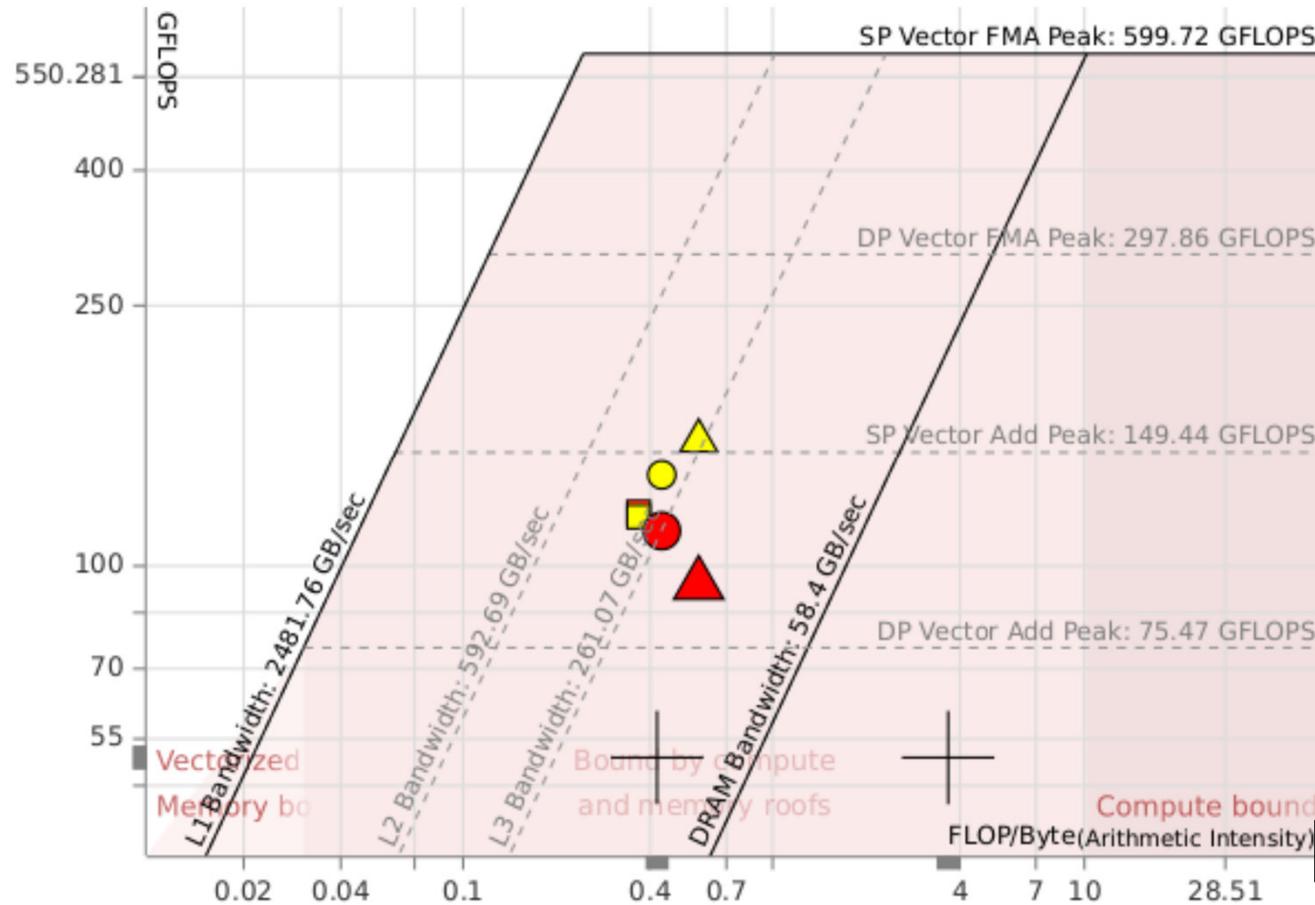
Intel® Xeon® Platinum
8180 (Skylake, 28 cores),
ICC v18.0, Devito v3.1

TTI (Tilted Transverse
Isotropy), second order in
time. 415 timesteps
(1000ms), single precision.

Space order:
4 (circles),
8 (crosses), and
12 (triangles)
16 (nablas)

Fabio Luporini et al. Architecture and Performance of Devito, a System for Automated Stencil Computation. ACM Trans. Math. Softw. 46, 1, Article 6 (April 2020), <https://doi.org/10.1145/3374916>

For latest performance data see
<https://www.devitoproject.org/thematrix/>



Single-socket 8-core Intel Broadwell E5-2673 v4 CPUs with AVX2, L1 (32KB), L2 (256KB) private to each core, 50MB shared L3 (Ubuntu 18.04.4, Devito v4.2.3)

Isotropic acoustic model, second-order in time, single-precision

Space order:

- 4 (triangles), ▲ ▲
- 8 (circles), and ● ●
- 12 (squares). ■ ■

Red markers show the performance of **spatially** blocked vectorized kernels

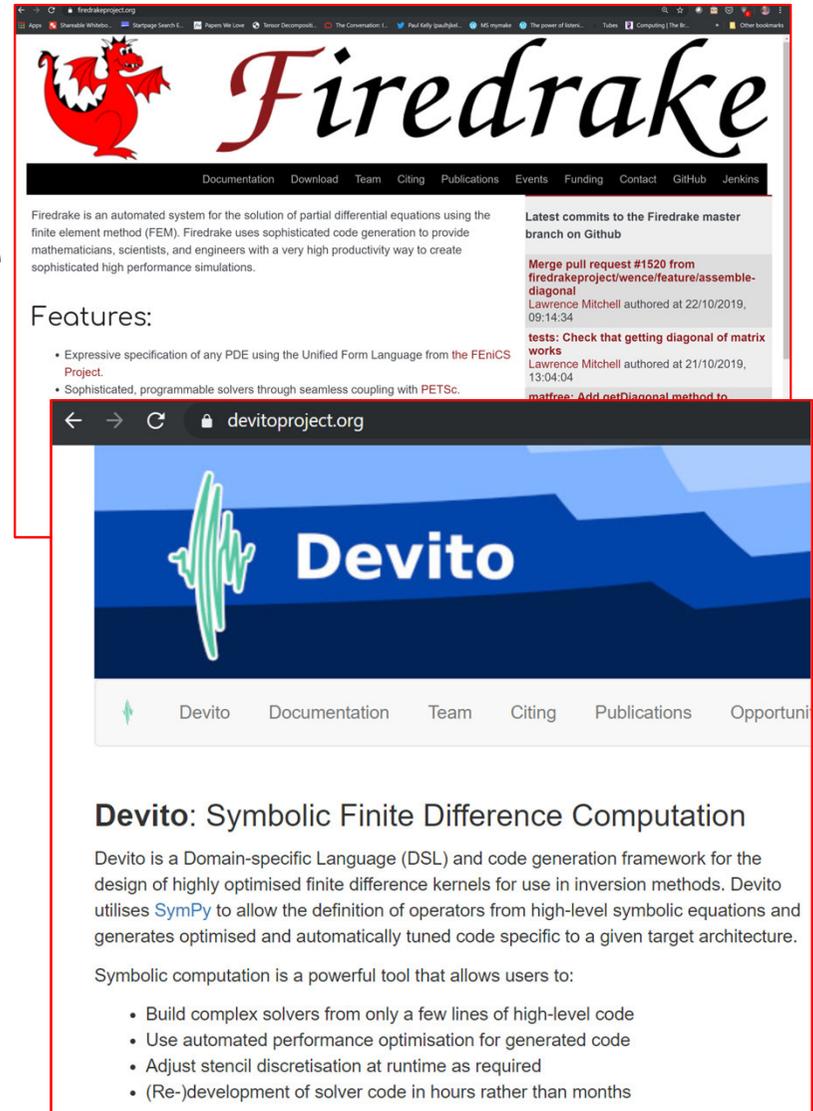
Yellow markers show **temporal** blocking using autotuned tile parameters.

George Bisbas, et al. **Temporal blocking of finite-difference stencil operators with sparse “off-the-grid” sources**. IPDPS 21 [arXiv:2010.10248](https://arxiv.org/abs/2010.10248)

- **Engaging with applications to exploit domain-specific optimisations can be incredibly fruitful**
 - Compiling general purpose languages is worthy but usually incremental
- **Compiler architecture is all about designing intermediate representations – that make hard things look easy**
 - Tools to deliver domain-specific optimisations often have domain-specific representations
 - Premature lowering is the constant enemy (appropriate lowering is great)
- **Along the way, we learn something about building better general-purpose compilers and programming abstractions**
 - Drill vertically, expand horizontally

How can we change the world?

- The real value of Firedrake and Devito is in supporting the applications users in exploring *their* design space
- We enable them to navigate rapidly through alternative solutions to their problem
- In the future, we will have automated pathways from maths to code for many classes of problem, and many alternative solution techniques

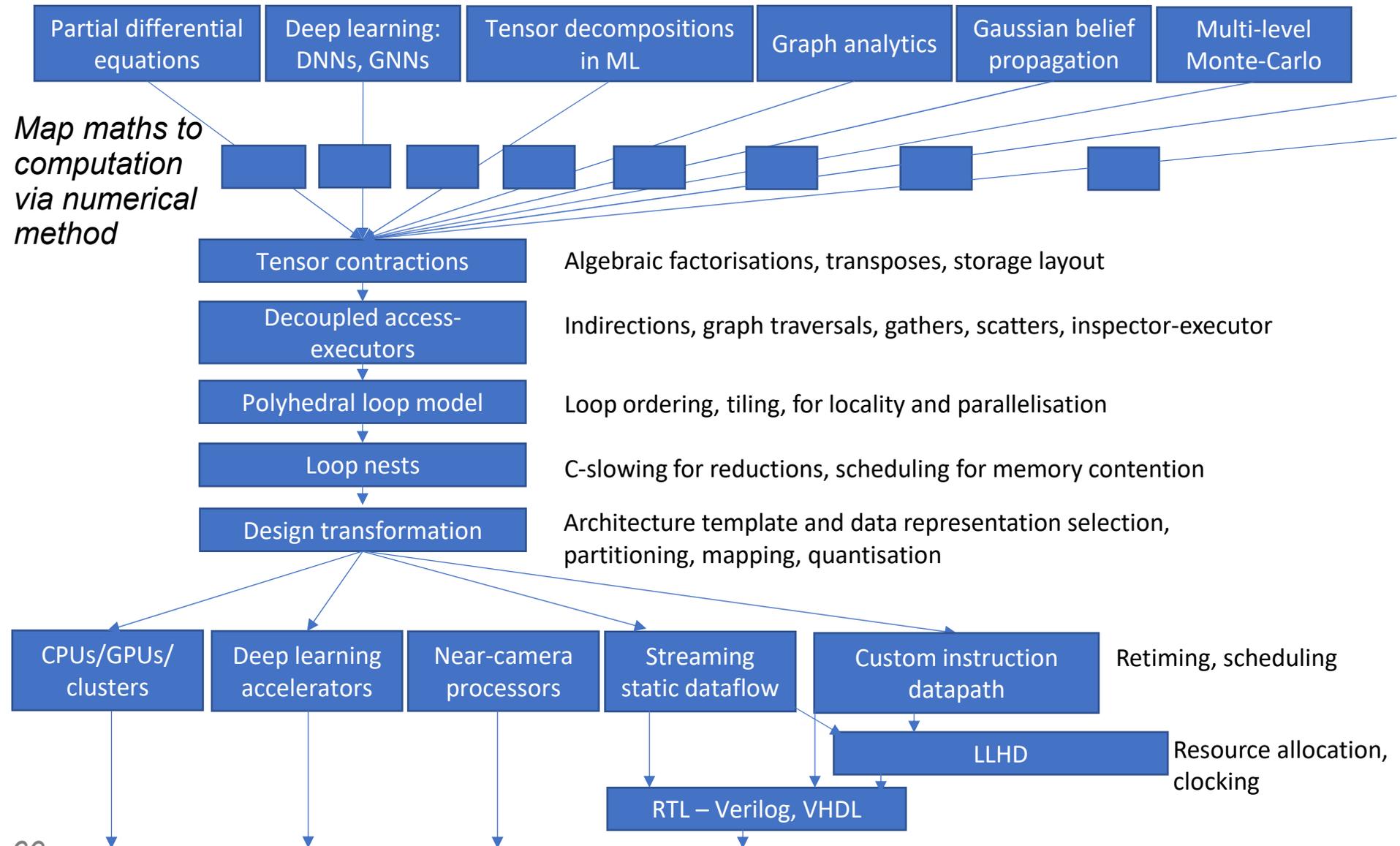


The image displays two screenshots of project websites. The top screenshot shows the Firedrake project website (firedrakeproject.org), featuring a red dragon logo and the word "Firedrake" in a large, stylized font. Below the header, there is a navigation menu with links for Documentation, Download, Team, Citing, Publications, Events, Funding, Contact, GitHub, and Jenkins. The main content area describes Firedrake as an automated system for solving partial differential equations using the finite element method (FEM). It lists features such as expressive specification of PDEs using the Unified Form Language (UFL) and sophisticated, programmable solvers through seamless coupling with PETSc. A sidebar on the right shows the latest commits to the master branch on GitHub, including a merge pull request and a test check.

The bottom screenshot shows the Devito project website (devitoproject.org). It features a blue header with a green waveform logo and the word "Devito" in a large, bold font. Below the header, there is a navigation menu with links for Devito, Documentation, Team, Citing, Publications, and Opportuni. The main content area describes Devito as a Domain-specific Language (DSL) and code generation framework for the design of highly optimised finite difference kernels. It mentions that Devito utilises SymPy to allow the definition of operators from high-level symbolic equations and generates optimised and automatically tuned code specific to a given target architecture. A section titled "Devito: Symbolic Finite Difference Computation" explains that symbolic computation is a powerful tool that allows users to:

- Build complex solvers from only a few lines of high-level code
- Use automated performance optimisation for generated code
- Adjust stencil discretisation at runtime as required
- (Re-)development of solver code in hours rather than months

Vision for the future



Thank you to our many many collaborators!

Partly funded/supported by

- SysGenX: Composable software generation for system-level simulation at Exascale (EP/W026066/1)
- XDSL: Efficient Cross-Domain DSL Development for Exascale (EP/W007789/1)
- NERC Doctoral Training Grant (NE/G523512/1)
- EPSRC “MAPDES” project (EP/I00677X/1)
- EPSRC “PSL” project (EP/I006761/1)
- Rolls Royce and the TSB through the SILOET programme
- EPSRC “PAMELA” Programme Grant (EP/K008730/1)
- EPSRC “PRISM” Platform Grants (EP/I006761/1 and EP/R029423/1)
- EPSRC “Custom Computing” Platform Grant (EP/I012036/1)
- EPSRC “Application Customisation” Platform Grant (EP/P010040/1)
- EPSRC “A new simulation and optimisation platform for marine technology” (EP/M011054/1)
- Basque Centre for Applied Mathematics (BCAM)

- Code:
 - <http://www.firedrakeproject.org/>
 - <http://op2.github.io/PyOP2/>
 - <https://github.com/OP-DSL/OP2-Common>

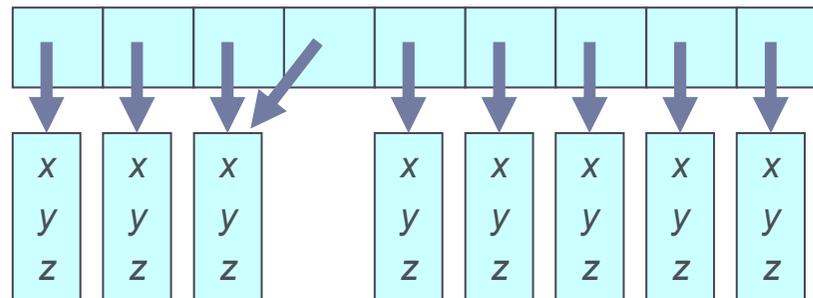
Extra slides for questions

Easy parallelism

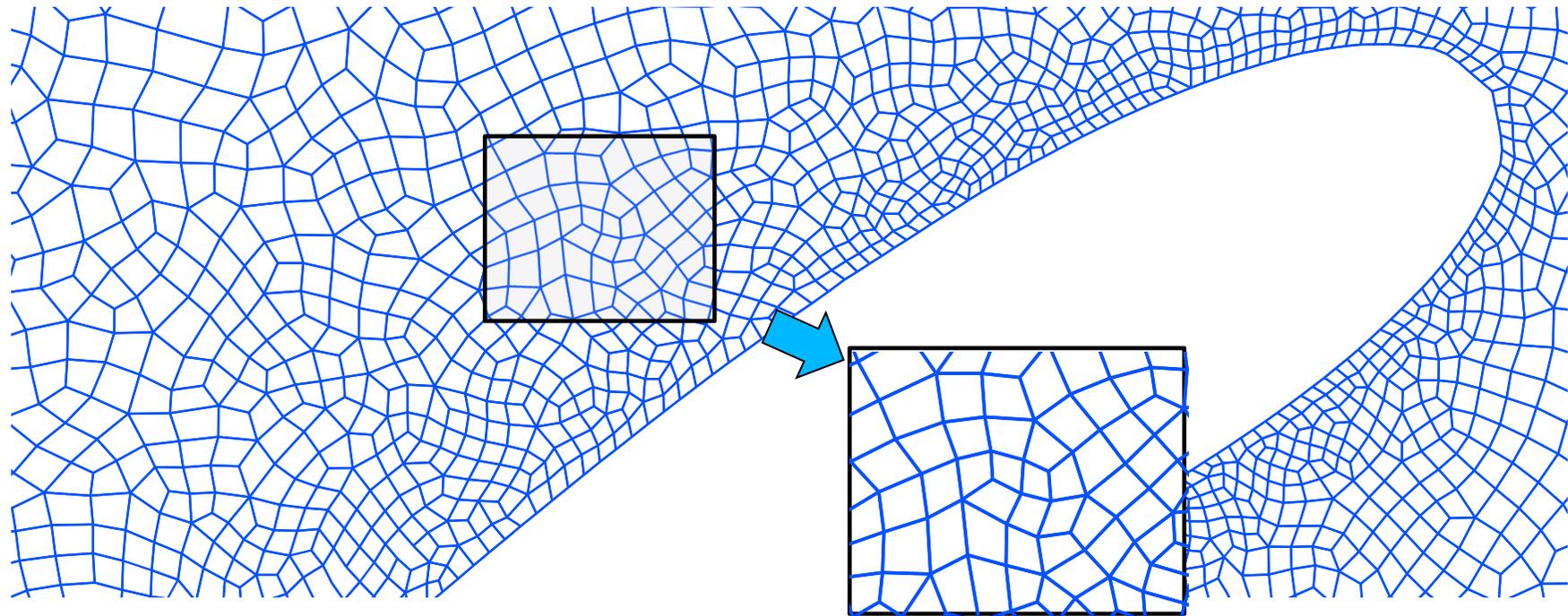
Example:

```
for (i=0; i<N; ++i) {  
    points[i]->x += 1;  
}
```

- Can the iterations of this loop be executed in parallel?

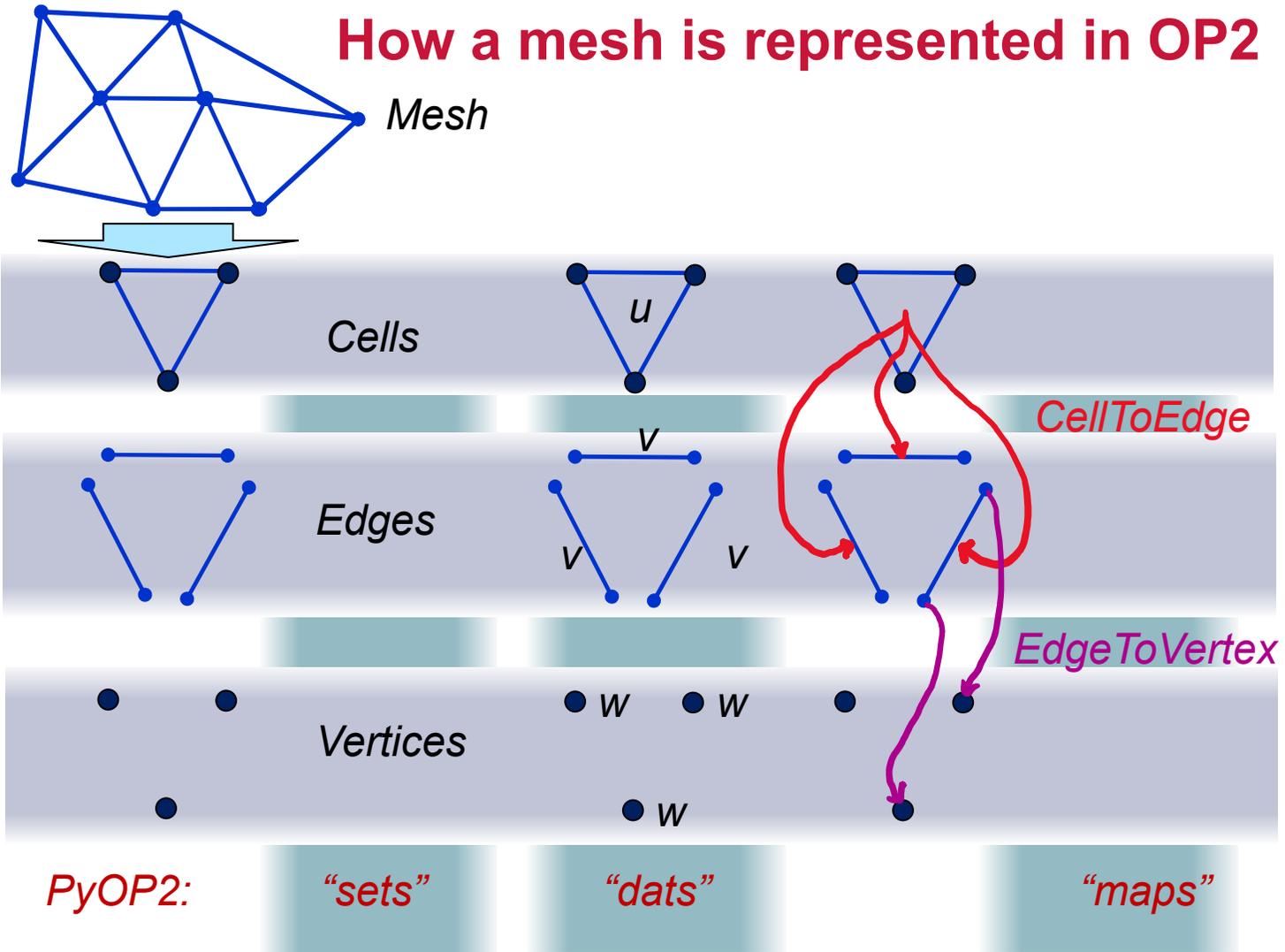


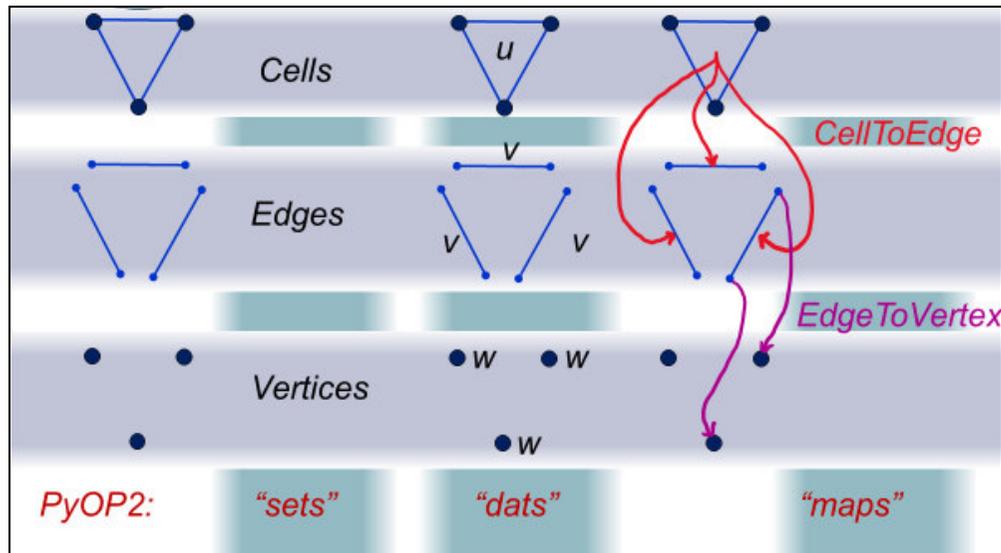
- Oh no: not all the iterations are independent!
 - You want to re-use piece of code in different contexts
 - Whether it's parallel depends on context!



- Unstructured meshes require pointers/indirection because adjacency lists have to be represented explicitly
- A controlled form of pointers (actually a general graph)
- **OP2** is a C++ and Fortran library for parallel loops over the mesh, implemented by source-to-source transformation
- **PyOP2** is the same basic model, implemented in Python using runtime code generation
- Enables generation of highly-optimised vectorised, CUDA, OpenMP and MPI code
- The OP2 model originates from Oxford (Mike Giles et al)

How a mesh is represented in OP2





OP2 loops, access descriptors and kernels

`op_par_loop(set, kernel, access descriptors)`

We specify which **set** to iterate over

We specify a **kernel** to execute – the kernel operates entirely locally, on the **dats** to which it has access

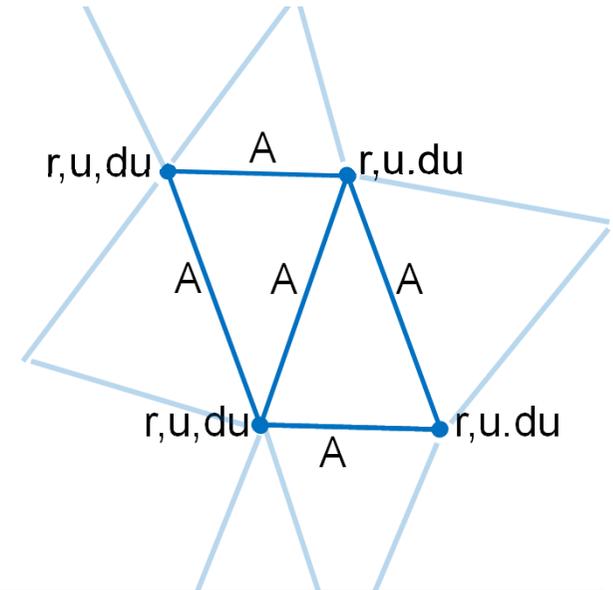
The **access descriptors** specify which dats the kernel has access to:

- **Which dats of the target set**
- **Which dats of sets indexed from this set through specified maps**

- OP2 separates local (kernel) from global (mesh)
- OP2 makes data dependence explicit

PyOP2: “decoupled access-execute”

- Parallel loops, over sets (nodes, edges etc)
- Access descriptors specify how to pass data to and from the C kernel
- The kernel operates only on local data



Access descriptors specify how to feed the kernel from the mesh

```
for iter in xrange(0, NITER):
```

```
    u_sum = op2.Global(1, data=0.0, np.float32)
```

```
    u_max = op2.Global(1, data=0.0, np.float32)
```

```
    op2.par_loop(res, edges,
                 p_A(op2.READ),
                 p_u(op2.READ, edge2vertex[1]),
                 p_du(op2.INC, edge2vertex[0]),
                 beta(op2.READ))
```

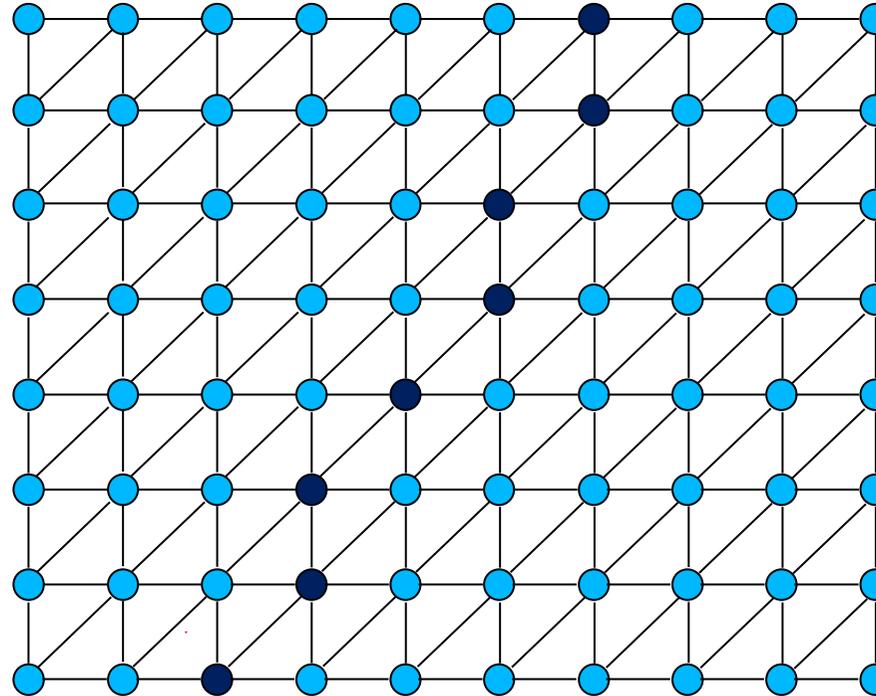
```
void res(float *A, float *u, float *du,
         const float *beta) {
    *du += (*beta) * (*A) * (*u);
}
```

```
    op2.par_loop(update, nodes,
                 p_r(op2.READ),
                 p_du(op2.RW),
                 p_u(op2.INC),
                 u_sum(op2.INC),
                 u_max(op2.MAX))
```

```
void update(float *r, float *du, float *u, float
            *u_sum, float *u_max) {
    *u += *du + alpha * (*r);
    *du = 0.0f;
    *u_sum += (*u) * (*u);
    *u_max = *u_max > *u ? *u_max : *u;
}
```

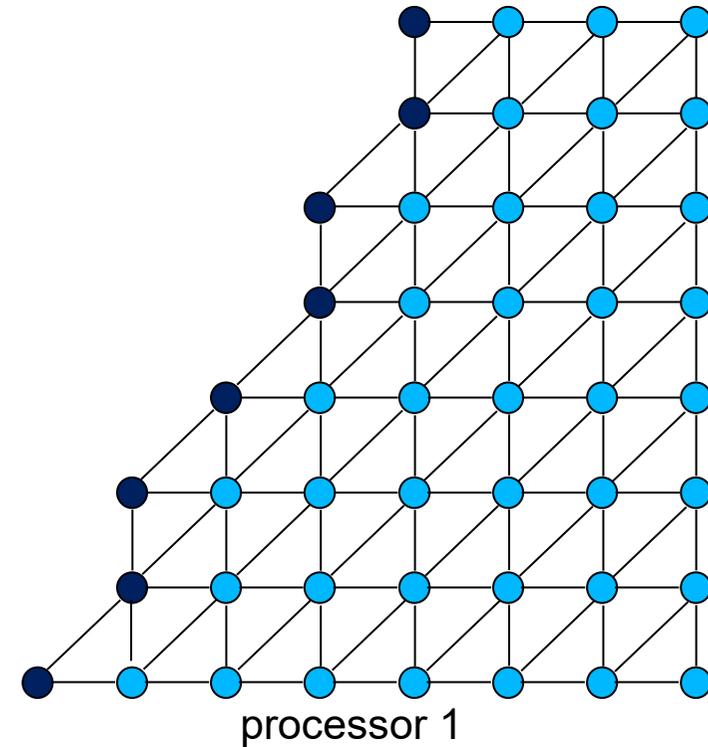
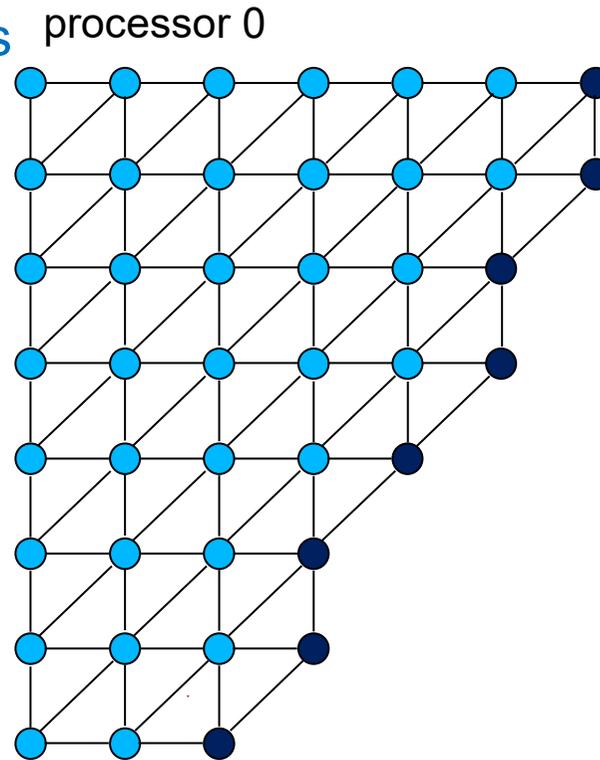
Code generation for indirect loops in PyOP2

- For MPI we precompute partitions & haloes
- Derived from PyOP2 access descriptors, implemented using PetSC DMPLex
- At partition boundaries, the entities (vertices, edges, cells) form layered halo region



Code generation for indirect loops in PyOP2

- For MPI we precompute partitions & haloes
- Derived from PyOP2 access descriptors, implemented using PetSC DMPLex
- At partition boundaries, the entities (vertices, edges, cells) form layered halo region



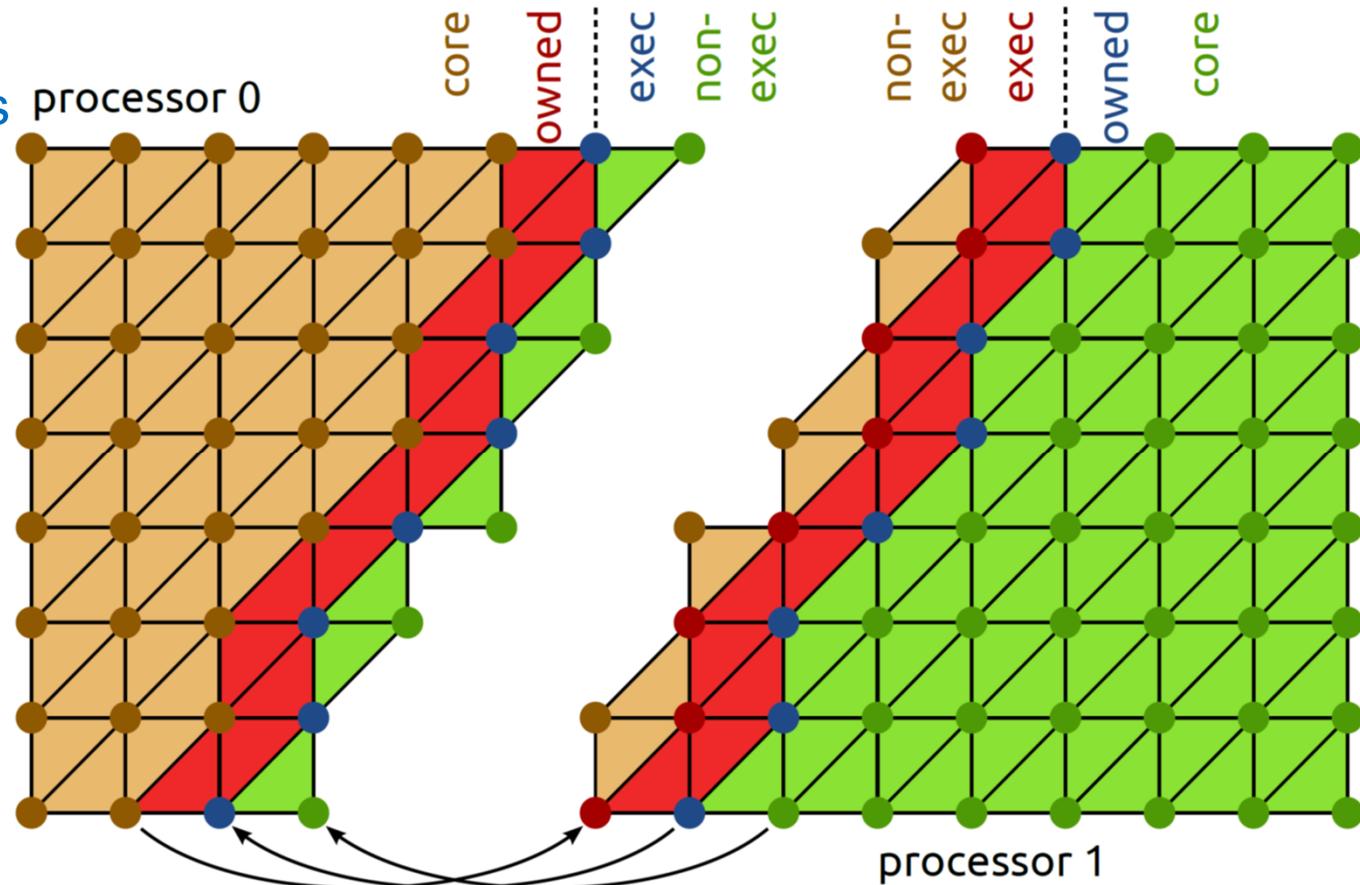
Code generation for indirect loops in PyOP2

- For MPI we precompute partitions & haloes

- Derived from PyOP2 access descriptors, implemented using PetSC DMPlex

- At partition boundaries, the entities (vertices, edges, cells) form layered halo region

- Core:** entities owned which can be processed without accessing halo data.
- Owned:** entities owned which access halo data when processed
- Exec halo:** off-processor entities which are redundantly executed over because they touch owned entities
- Non-exec halo:** off-processor entities which are not processed, but read when computing the exec halo



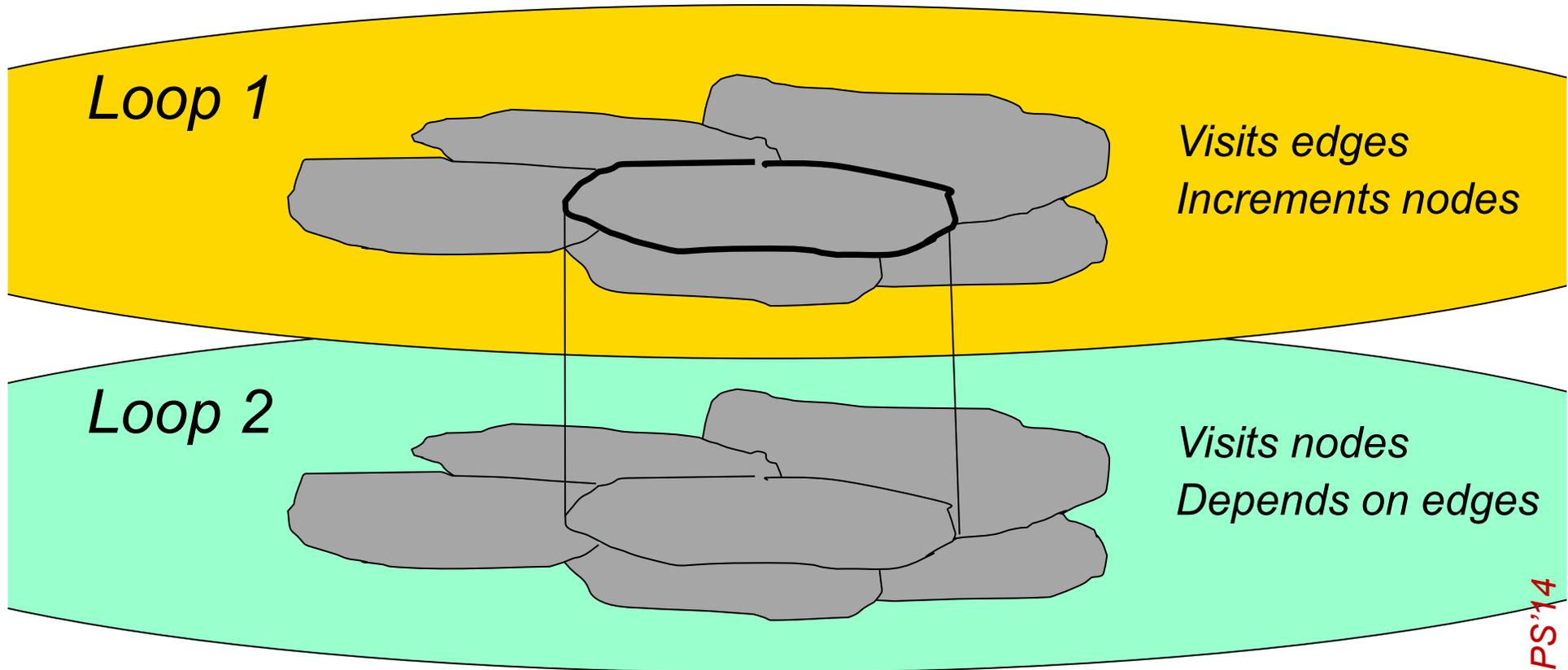
■ Can we automate interesting optimisations that would be hard to do by hand?

■ First example:

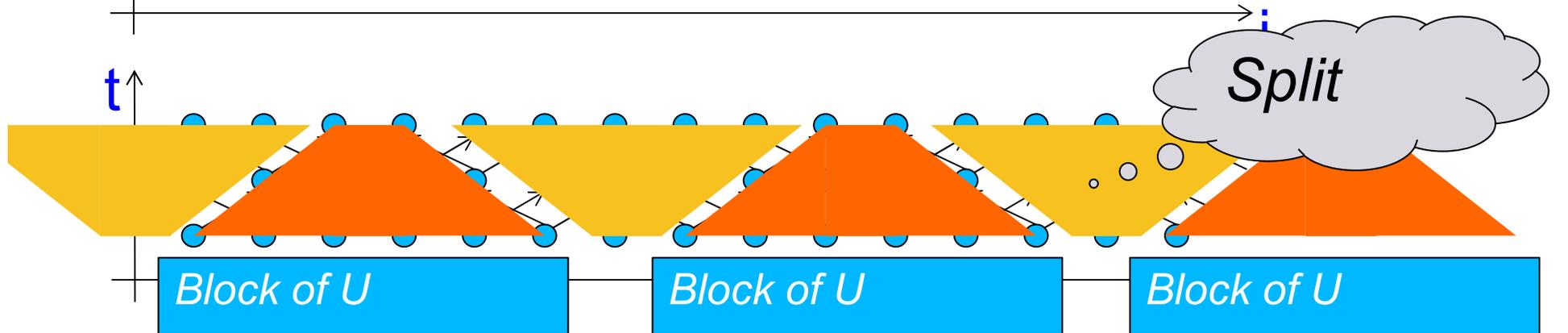
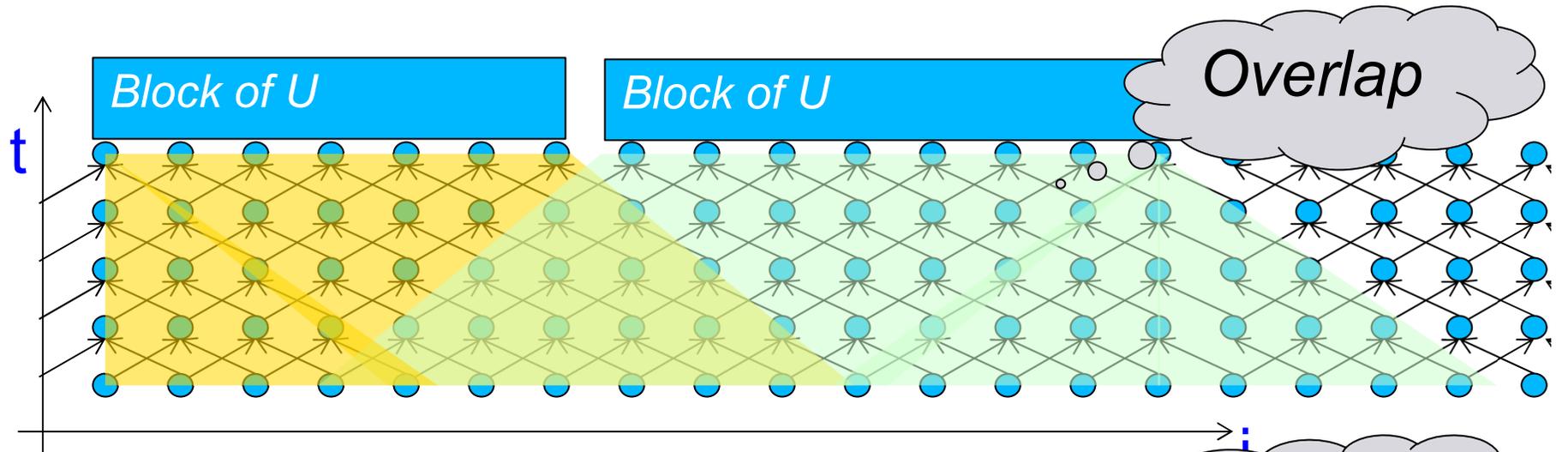
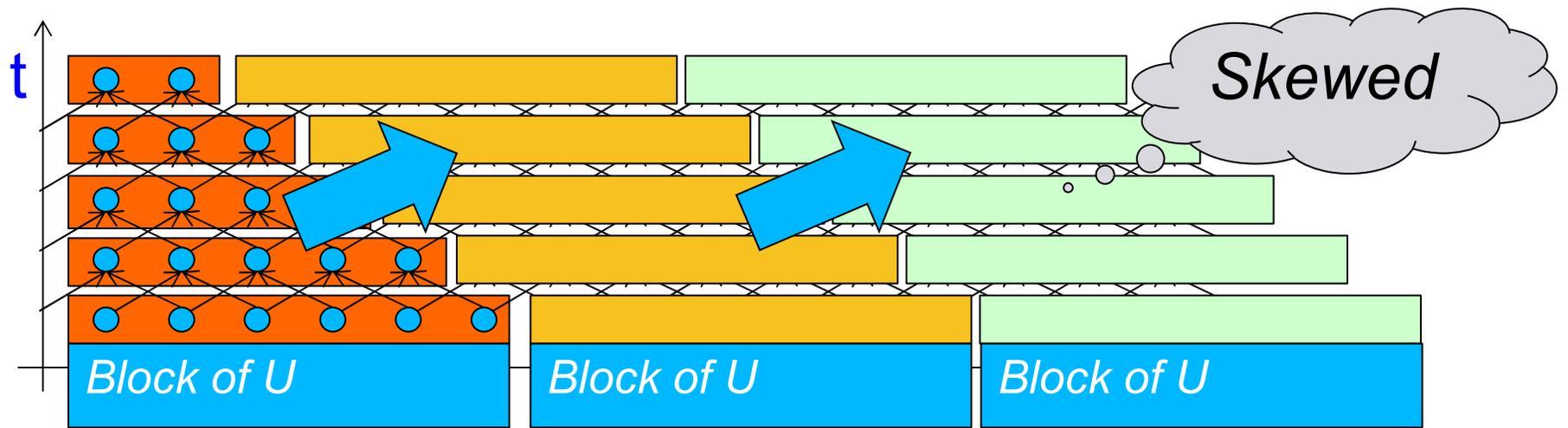
■ Tiling for cache locality

■ (This optimisation has been implemented – and automated – but does not currently form part of the standard distribution)

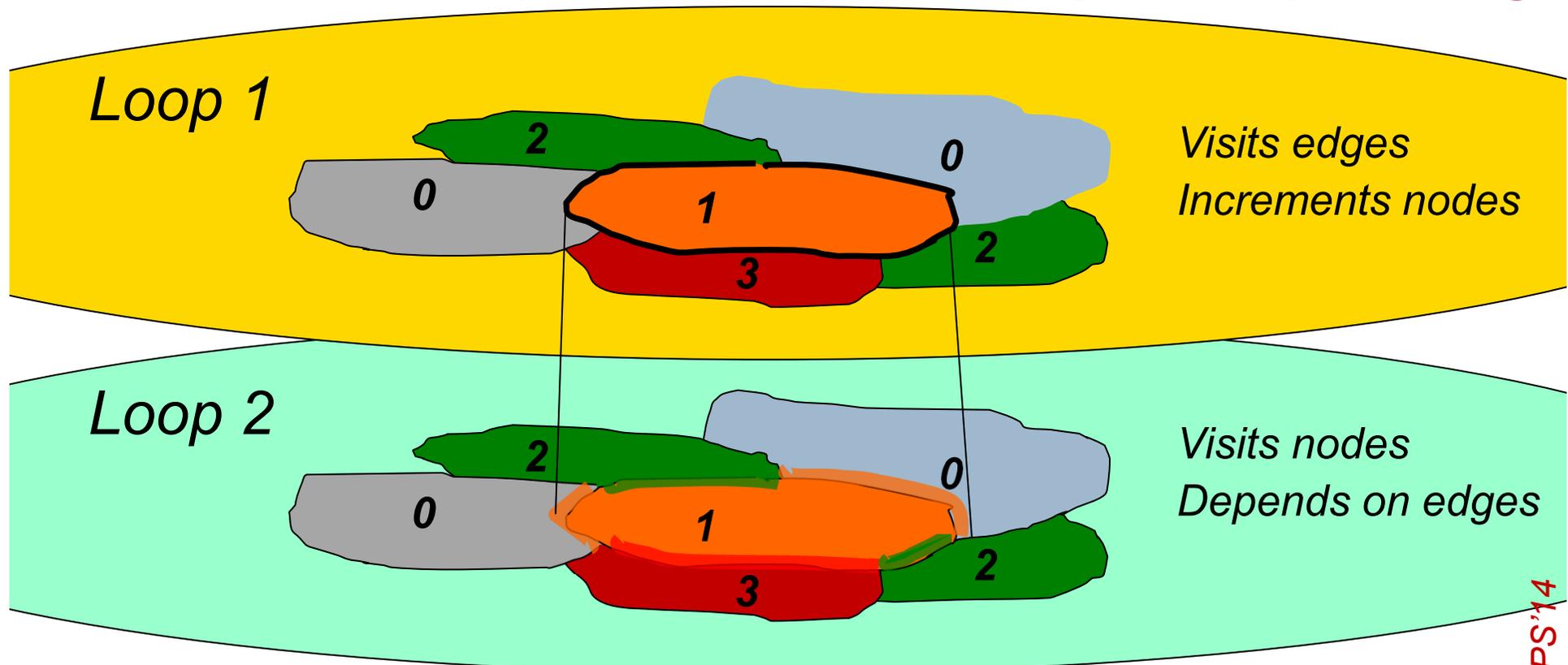
Sparse split tiling on an unstructured mesh, for locality



- How can we load a block of mesh and do the iterations of loop 1, then the iterations of loop 2, before moving to the next block?
- If we could, we could dramatically improve the memory access behaviour!

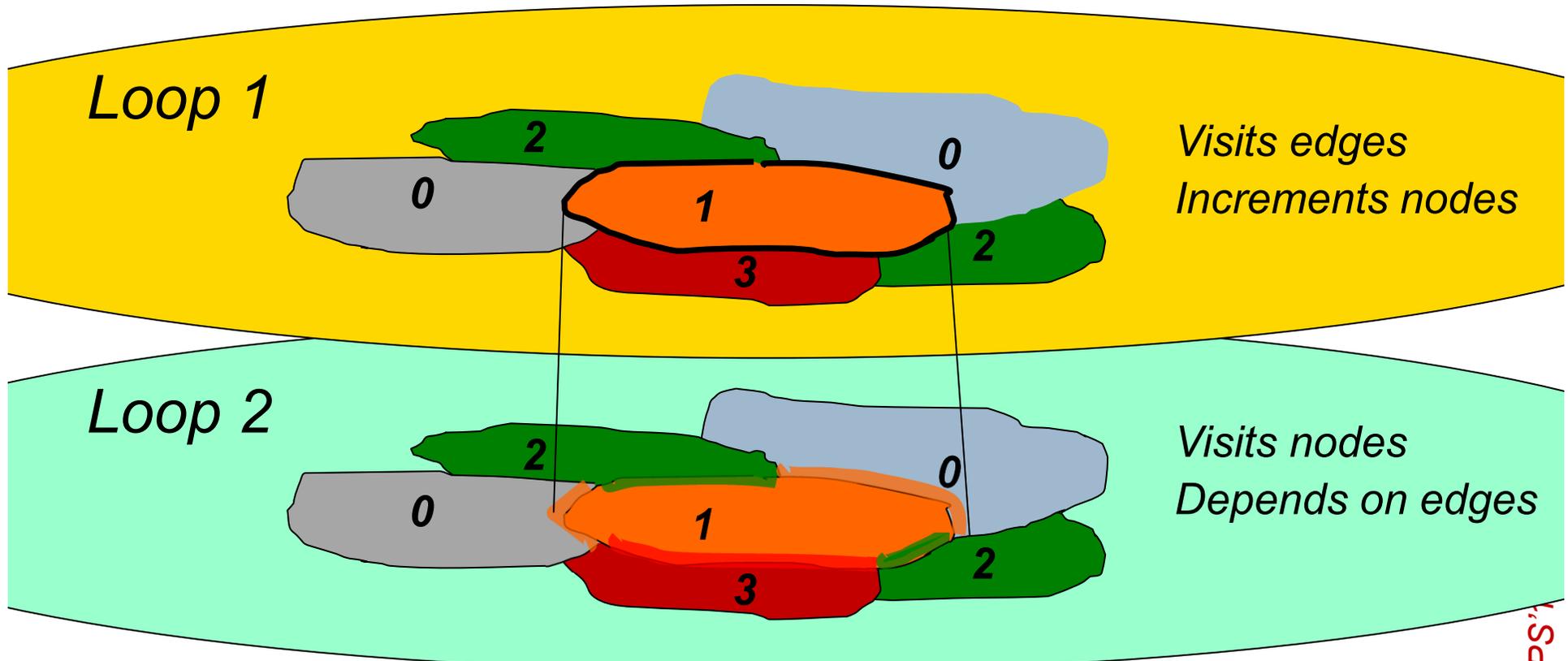


Sparse split tiling



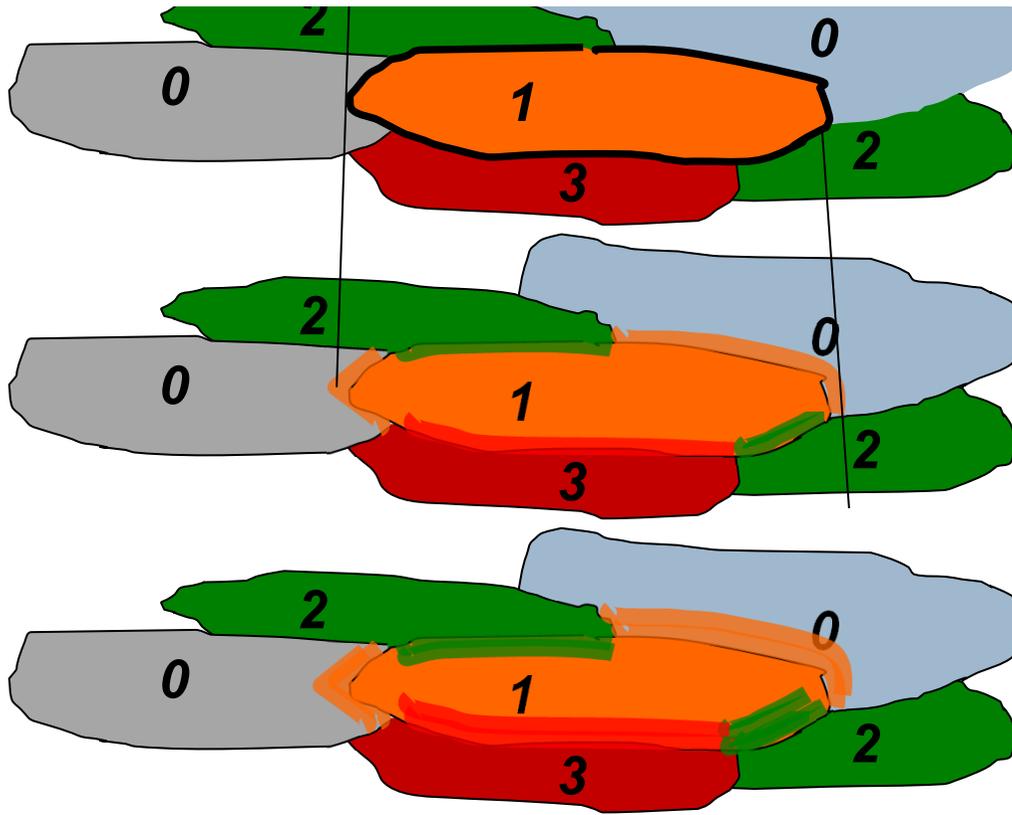
- Partition the iteration space of loop 1
- Colour the partitions, execute the colours in order
- Project the tiles, using the knowledge that colour n can use data produced by colour $n-1$
- Thus, the tile coloured #1 **grows** where it meets colour #0
- And **shrinks** where it meets colours #2 and #3

Sparse split tiling

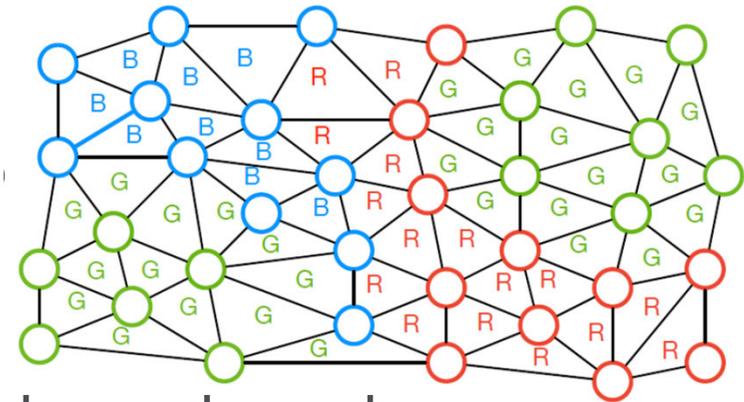
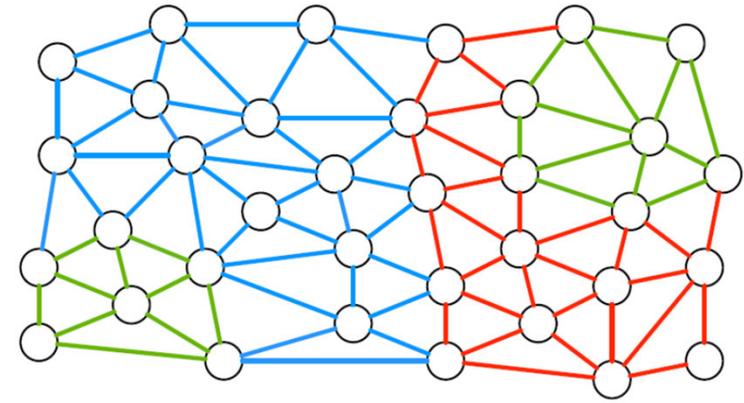


- Partition the iteration space of loop 1
- Colour the partitions
- Project the tiles, using the knowledge data produced by colour n-1
- Thus, the tile coloured #1 grows when
- And shrinks where it meets colours #2 and #3

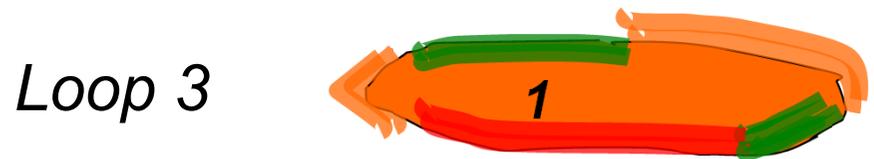
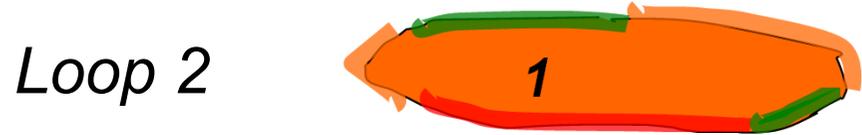
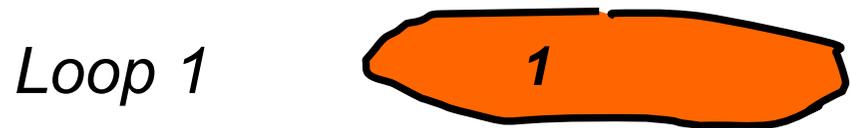
*Inspector-executor:
derive tasks and
task graph from
the mesh, **at
runtime***



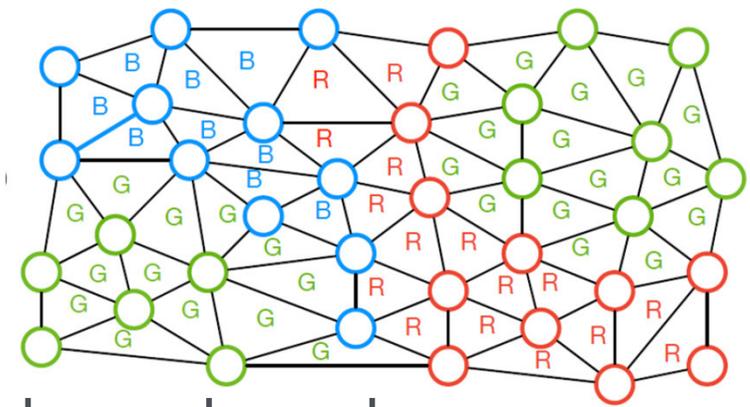
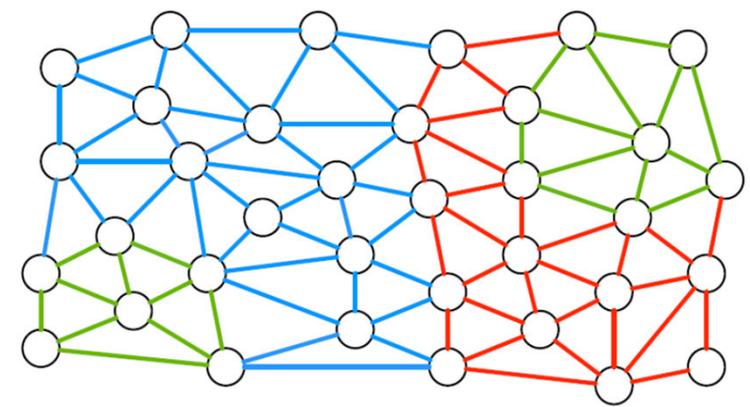
Tiles grow



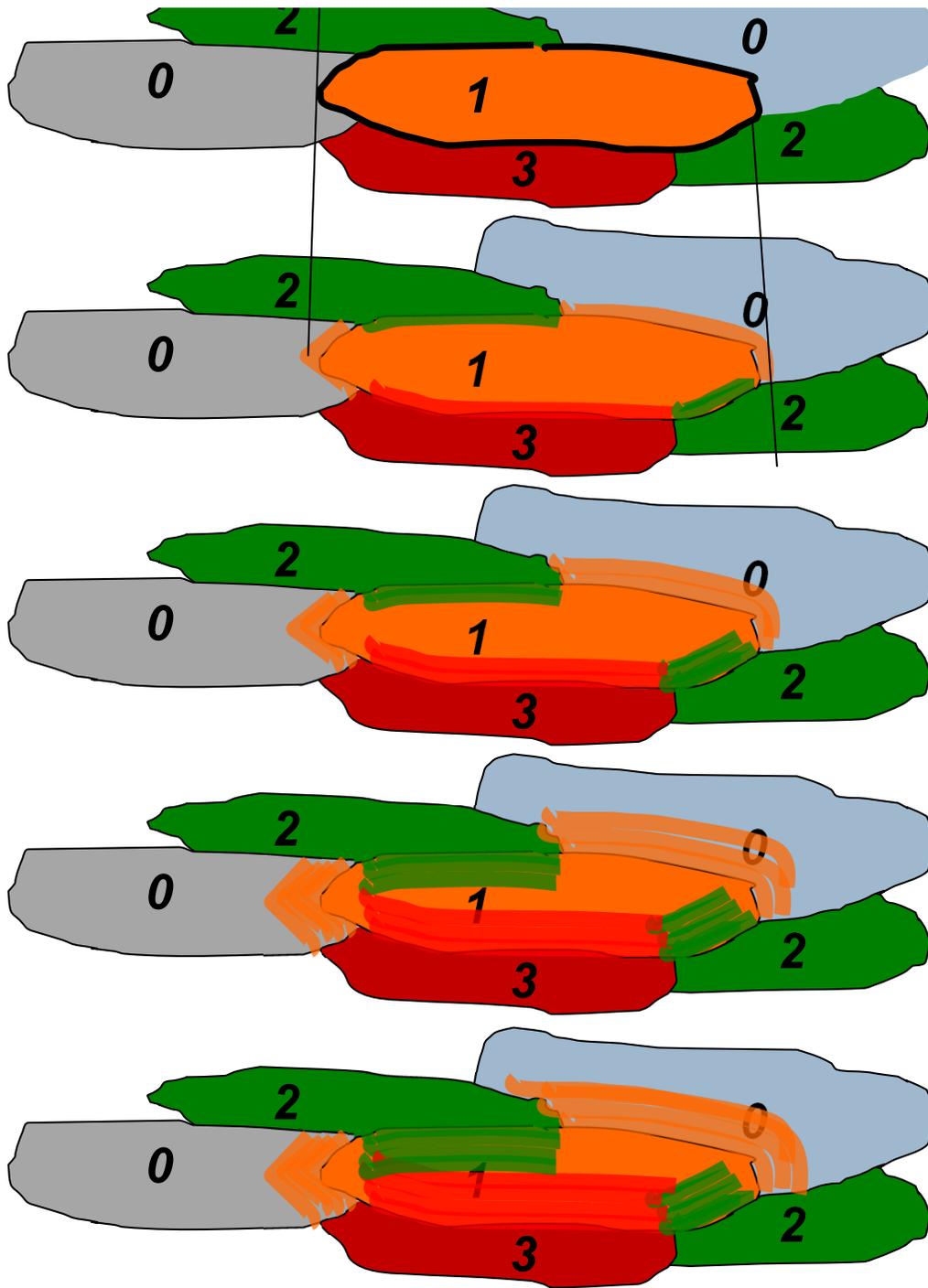
- As we project the tiles forward, tile shape degrades
- Perimeter-volume ratio gets worse



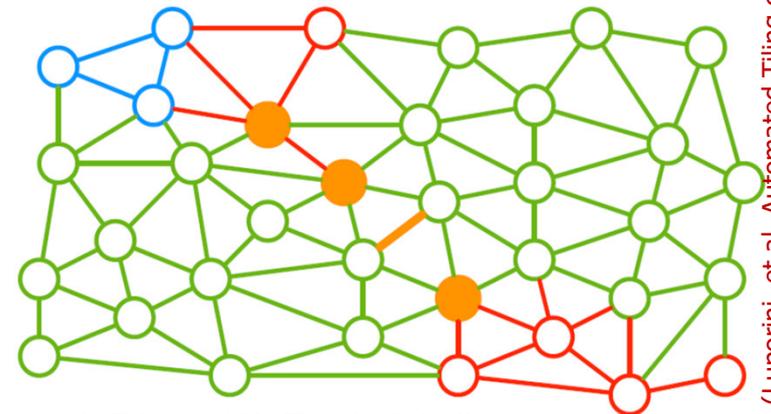
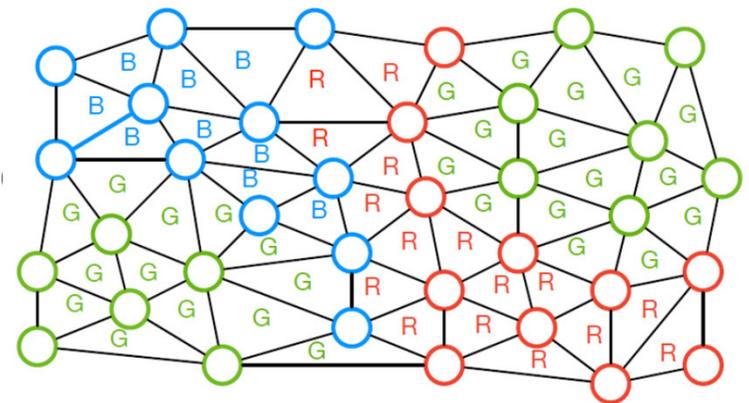
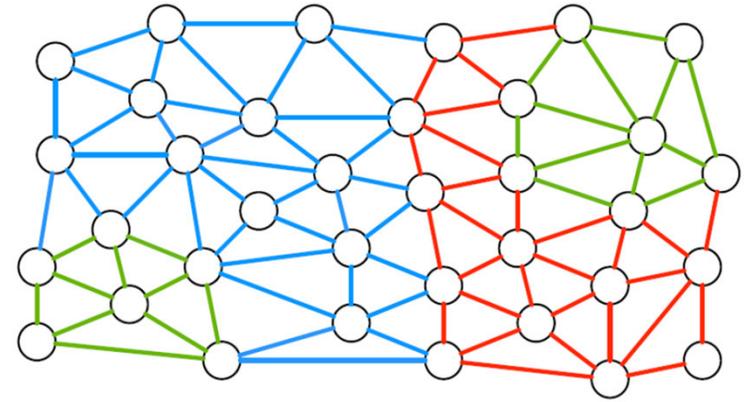
Tiles grow



- As we project the tiles forward, tile shape degrades
- Perimeter-volume ratio gets worse
- We could partition Loop 1's data for the cache
- But Loop 2 and Loop 3 have different footprints
- So we rely on good (ideally space-filling-curve) numbering



Tiles can collide

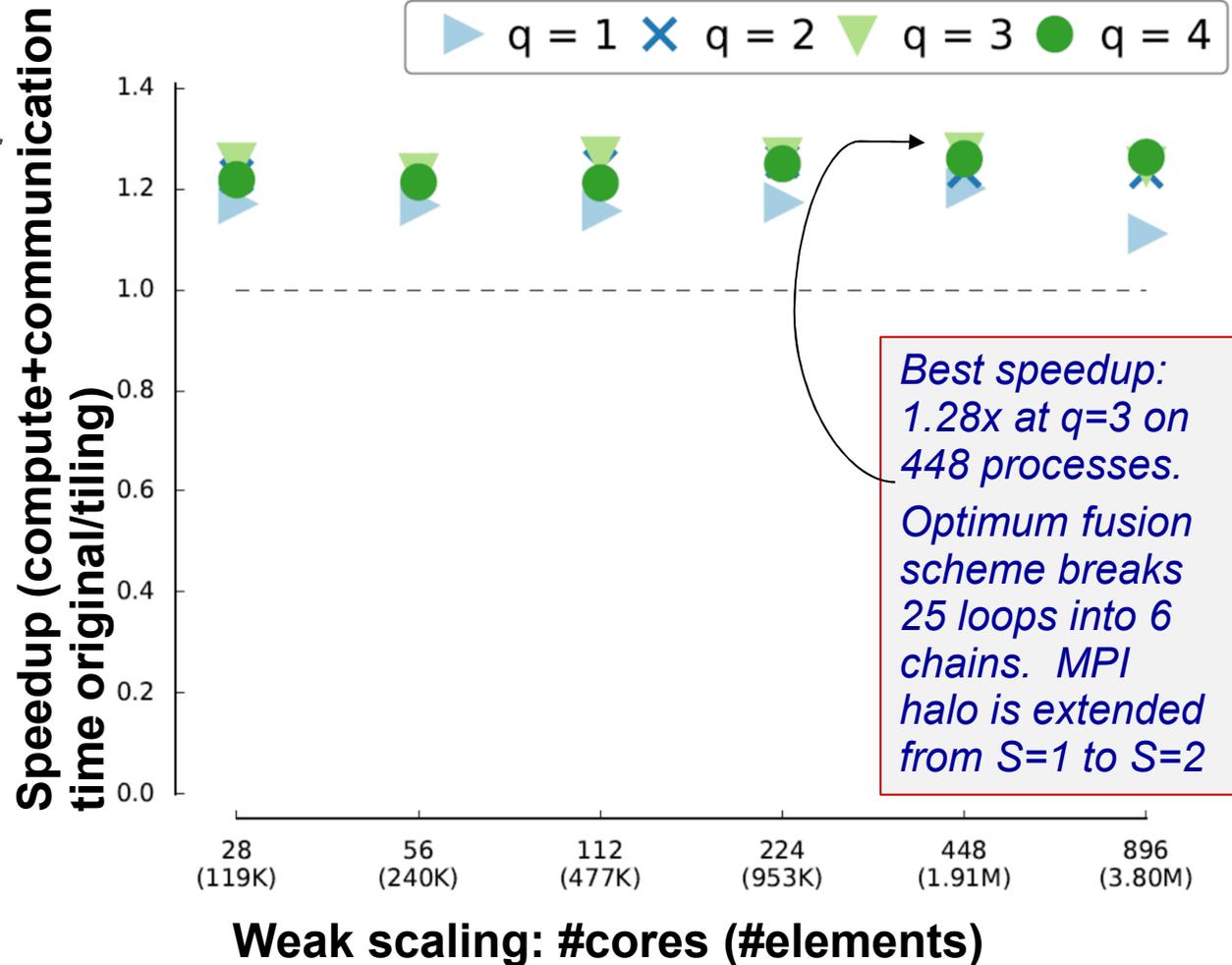


(1) Blue, (2) Red, (3) Green

(Luporini, et al, Automated Tiling of Unstructured Mesh Computations with Application to Seismological Modeling. ACM TOMS 2019)

■ Example: Seigen

- Elastic wave solver
- 2d triangular mesh
- Velocity-stress formulation
- 4th-order explicit leapfrog timestepping scheme
- Discontinuous-Galerkin, order $q=1-4$
- 32 nodes, 2x14-core E5-2680v4, SGI MPT 2.14
- 1000 timesteps (ca. 1.15s/timestep)



- Up to 1.28x speedup
- Inspection about as much time as 2 timesteps
- Using RCM numbering – space-filling curve should lead to better results

■ Can we automate interesting optimisations that would be hard to do by hand?

■ Second example:

■ Generalised loop-invariant code motion

■ (This optimisation has been implemented, automated, and re-implemented – and forms part of the standard distribution)

Recall:

```
double form_t2;
double const form_t20[7 * 6] = { ... };
double form_t21...t37;
double form_t38[6];
double form_t39[6];
double form_t4;
double form_t40...t45;
double form_t5...t9;
double t0[6 * 2];
double t1[3 * 2];
double t2[6 * 2 * 6 * 2];

for (int n = start; n <= -1 + end; ++n)
{
  for (int i4 = 0; i4 <= 5; ++i4)
    for (int i5 = 0; i5 <= 1; ++i5)
      for (int i6 = 0; i6 <= 5; ++i6)
        for (int i7 = 0; i7 <= 1; ++i7)
          t2[24 * i4 + i2 * 15 + 2 * i6 + i7] = 0.0;
  for (int i2 = 0; i2 <= 2; ++i2)
    for (int i3 = 0; i3 <= 1; ++i3)
      t1[2 * i2 + i3] = dat1[2 * map1[3 * n + i2] + i3];
  for (int i0 = 0; i0 <= 5; ++i0)
    for (int i1 = 0; i1 <= 1; ++i1)
      t0[2 * i0 + i1] = dat0[2 * map0[6 * n + i0] + i1];
  form_t0 = -1.0 * t1[i1];
  form_t1 = form_t0 + t1[3];
  form_t2 = -1.0 * t1[0];
  form_t3 = form_t2 + t1[2];
  form_t4 = form_t0 + t1[5];
  form_t5 = form_t2 + t1[4];
  form_t6 = form_t3 * form_t4 + -1.0 * form_t5 * form_t1;
  form_t7 = 1.0 / form_t6;
  form_t8 = form_t7 * -1.0 * form_t1;
  form_t9 = form_t4 * form_t7;
  form_t10 = form_t3 * form_t7;
  form_t11 = form_t7 * -1.0 * form_t5;
  form_t12 = 0.0001 * (form_t8 * form_t9 + form_t10 * form_t11);
  form_t13 = 0.0001 * (form_t8 * form_t8 + form_t10 * form_t10);
  form_t14 = 0.0001 * (form_t9 * form_t9 + form_t11 * form_t11);
  form_t15 = 0.0001 * (form_t9 * form_t8 + form_t11 * form_t10);
  form_t16 = fabs(form_t6);
  for (int form_ip = 0; form_ip <= 6; ++form_ip)
  {
    form_t26 = 0.0; form_t25 = 0.0; form_t24 = 0.0; form_t23 = 0.0; form_t22 = 0.0; form_t21 = 0.0;
    for (int form_i = 0; form_i <= 5; ++form_i)
    {
      form_t21 = form_t21 + form_t20[6 * form_ip + form_i] * t0[1 + 2 * form_i];
      form_t22 = form_t22 + form_t19[6 * form_ip + form_i] * t0[1 + 2 * form_i];
      form_t23 = form_t23 + form_t20[6 * form_ip + form_i] * t0[2 * form_i];
      form_t24 = form_t24 + form_t19[6 * form_ip + form_i] * t0[2 * form_i];
      form_t25 = form_t25 + form_t18[6 * form_ip + form_i] * t0[1 + 2 * form_i];
      form_t26 = form_t26 + form_t18[6 * form_ip + form_i] * t0[2 * form_i];
    }
    form_t27 = form_t17[form_ip] * form_t16;
    form_t28 = form_t27 * form_t15;
    form_t29 = form_t27 * form_t14;
    form_t30 = form_t27 * (form_t26 * form_t9 + form_t25 * form_t11);
    form_t31 = form_t27 * form_t13;
    form_t32 = form_t27 * form_t12;
    form_t33 = form_t27 * (form_t26 * form_t8 + form_t25 * form_t10);
    form_t34 = form_t27 * (form_t11 * form_t24 + form_t10 * form_t23);
    form_t35 = form_t27 * (form_t9 * form_t22 + form_t8 * form_t21);
    form_t36 = form_t27 * (50.0 + form_t9 * form_t24 + form_t8 * form_t23);
    form_t37 = form_t27 * (50.0 + form_t11 * form_t22 + form_t10 * form_t21);
    for (int form_k0 = 0; form_k0 <= 5; ++form_k0)
    {
      form_t38[form_k0] = form_t18[6 * form_ip + form_k0] * form_t37;
      form_t39[form_k0] = form_t18[6 * form_ip + form_k0] * form_t36;
    }
    for (int form_j0 = 0; form_j0 <= 5; ++form_j0)
    {
      form_t40 = form_t18[6 * form_ip + form_j0] * form_t35;
      form_t41 = form_t18[6 * form_ip + form_j0] * form_t34;
      form_t42 = form_t20[6 * form_ip + form_j0] * form_t31 + form_t18[6 * form_ip + form_j0] * form_t33 + form_t19[6 * form_ip + form_j0] * form_t32;
      form_t43 = form_t20[6 * form_ip + form_j0] * form_t28 + form_t18[6 * form_ip + form_j0] * form_t30 + form_t19[6 * form_ip + form_j0] * form_t29;
      for (int form_k0_0 = 0; form_k0_0 <= 5; ++form_k0_0)
      {
        form_t44 = form_t43 * form_t19[6 * form_ip + form_k0_0];
        form_t45 = form_t42 * form_t20[6 * form_ip + form_k0_0];
        t2[24 * form_j0 + 2 * form_k0_0] = t2[24 * form_j0 + 2 * form_k0_0] + form_t45 + form_t18[6 * form_ip + form_j0] * form_t39[form_k0_0] + form_t44;
        t2[13 + 24 * form_j0 + 2 * form_k0_0] = t2[13 + 24 * form_j0 + 2 * form_k0_0] + form_t13 + 24 * form_j0 + 2 * form_k0_0] + form_t18[6 * form_ip + form_j0] * form_t41;
        t2[1 + 24 * form_j0 + 2 * form_k0_0] = t2[1 + 24 * form_j0 + 2 * form_k0_0] + form_t18[6 * form_ip + form_k0_0] * form_t41;
        t2[12 + 24 * form_j0 + 2 * form_k0_0] = t2[12 + 24 * form_j0 + 2 * form_k0_0] + form_t18[6 * form_ip + form_k0_0] * form_t40;
      }
    }
  }
  MatSetValuesBlockedLocal(mat0, 6, &(map0[6 * n]), 6, &(map0[6 * n]), &(t2[0]), ADD_VALUES);
}
```

```
t _dat1, double const *_restrict_ dat0, int const *_restrict_ map0, int const *_restrict_ map1)
```

- Generated code to assemble the resulting linear system matrix
- Executed at each triangle in the mesh
- Accesses degrees of freedom shared with neighbour triangles through indirection map

A simpler example:

```
void helmholtz(double A[3][3], double **coords) {  
  // K, det = Compute Jacobian (coords)
```

```
  static const double W[3] = {...}  
  static const double X_D10[3][3] = {{...}}  
  static const double X_D01[3][3] = {{...}}
```

```
  for (int i = 0; i < 3; i++)  
    for (int j = 0; j < 3; j++)  
      for (int k = 0; k < 3; k++)  
        A[j][k] += ((Y[i][k]*Y[i][j]+  
          +((K1*X_D10[i][k]+K3*X_D01[i][k])*(K1*X_D10[i][j]+K3*X_D01[i][j]))+  
          +((K0*X_D10[i][k]+K2*X_D01[i][k])*(K0*X_D10[i][j]+K2*X_D01[i][j]))))*  
          *det*W[i]);  
}
```

- Local assembly code generated by Firedrake for a Helmholtz problem on a 2D triangular mesh using Lagrange $p = 1$ elements.
- The local assembly operation computes a small dense submatrix
- These are combined to form a global system of simultaneous equations capturing the discretised conservation laws expressed by the PDE

A simpler example:

```
void helmholtz(double A[3][3], double **coords) {  
    // K, det = Compute Jacobian (coords)
```

```
    static const double W[3] = {...}  
    static const double X_D10[3][3] = {{...}}  
    static const double X_D01[3][3] = {{...}}
```

```
    for (int i = 0; i < 3; i++)  
        for (int j = 0; j < 3; j++)  
            for (int k = 0; k < 3; k++)  
                A[j][k] += ((Y[i][k]*Y[i][j]+  
                    +((K1*X_D10[i][k]+K3*X_D01[i][k])*(K1*X_D10[i][j]+K3*X_D01[i][j]))+  
                    +((K0*X_D10[i][k]+K2*X_D01[i][k])*(K0*X_D10[i][j]+K2*X_D01[i][j])))  
                    *det*W[i]);  
    }
```

- Local assembly code generated by Firedrake for a Helmholtz problem on a 2D triangular mesh using Lagrange $p = 1$ elements.
- The local assembly operation computes a small dense submatrix
- These are combined to form a global system of simultaneous equations capturing the discretised conservation laws expressed by the PDE

Generalised loop-invariant code motion:

```
void helmholtz(double A[3][4], double **coords) {
  #define ALIGN __attribute__((aligned(32)))
  // K, det = Compute Jacobian (coords)

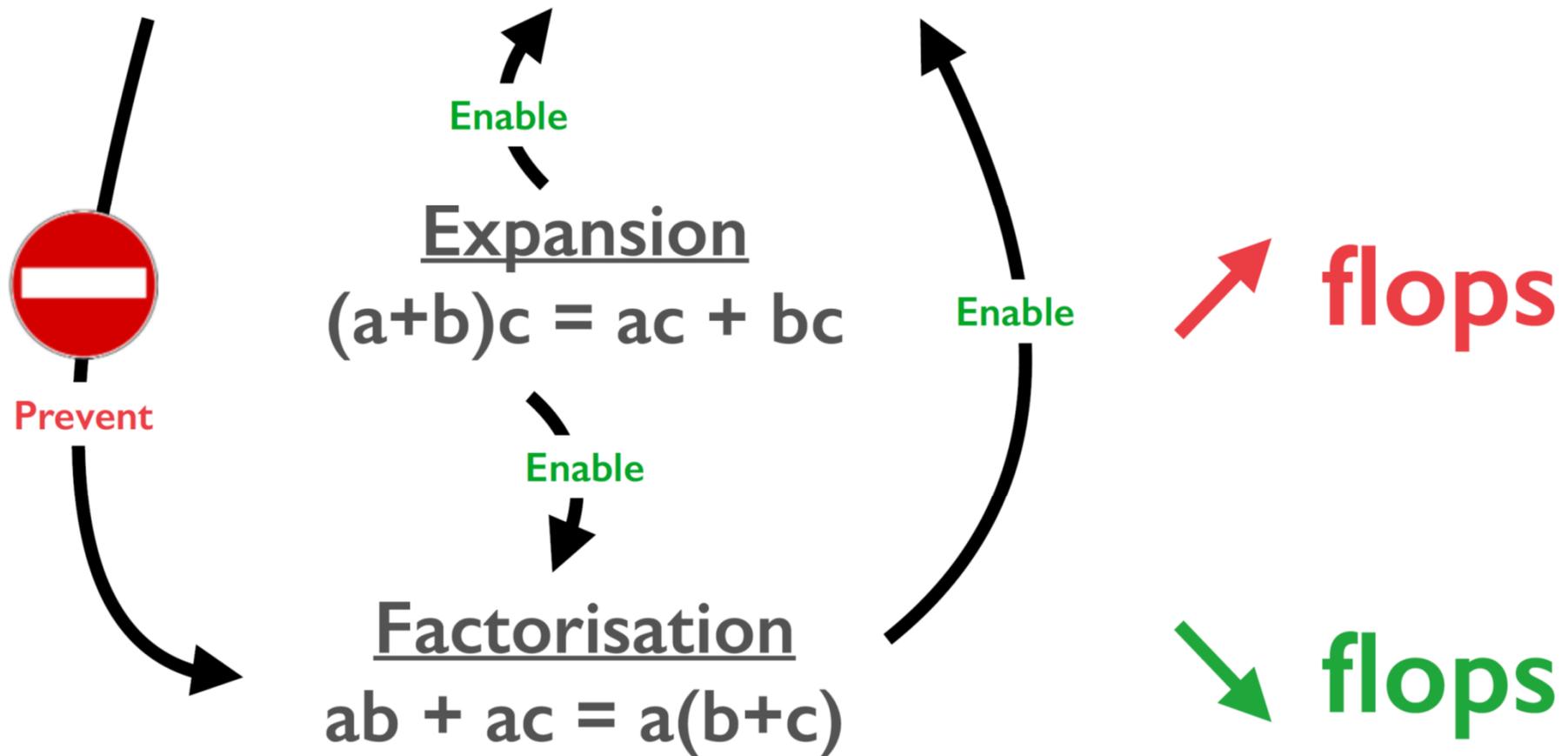
  static const double W[3] ALIGN = {...}
  static const double X_D10[3][4] ALIGN = {...}
  static const double X_D01[3][4] ALIGN = {...}

  for (int i = 0; i<3; i++) {
    double LI_0[4] ALIGN;
    double LI_1[4] ALIGN;
    for (int r = 0; r<4; r++) {
      LI_0[r] = ((K1*X_D10[i][r])+(K3*X_D01[i][r]));
      LI_1[r] = ((K0*X_D10[i][r])+(K2*X_D01[i][r]));
    }
    for (int j = 0; j<3; j++)
      #pragma vector aligned
      for (int k = 0; k<4; k++)
        A[j][k] += (Y[i][k]*Y[i][j]+LI_0[k]*LI_0[j]+LI_1[k]*LI_1[j])*det*W[i];
  }
}
```

- Local assembly code for the Helmholtz problem after application of
 - padding,
 - data alignment,
 - Loop-invariant code motion
- In this example, sub-expressions invariant to j are identical to those invariant to k , so they can be precomputed once in the r loop

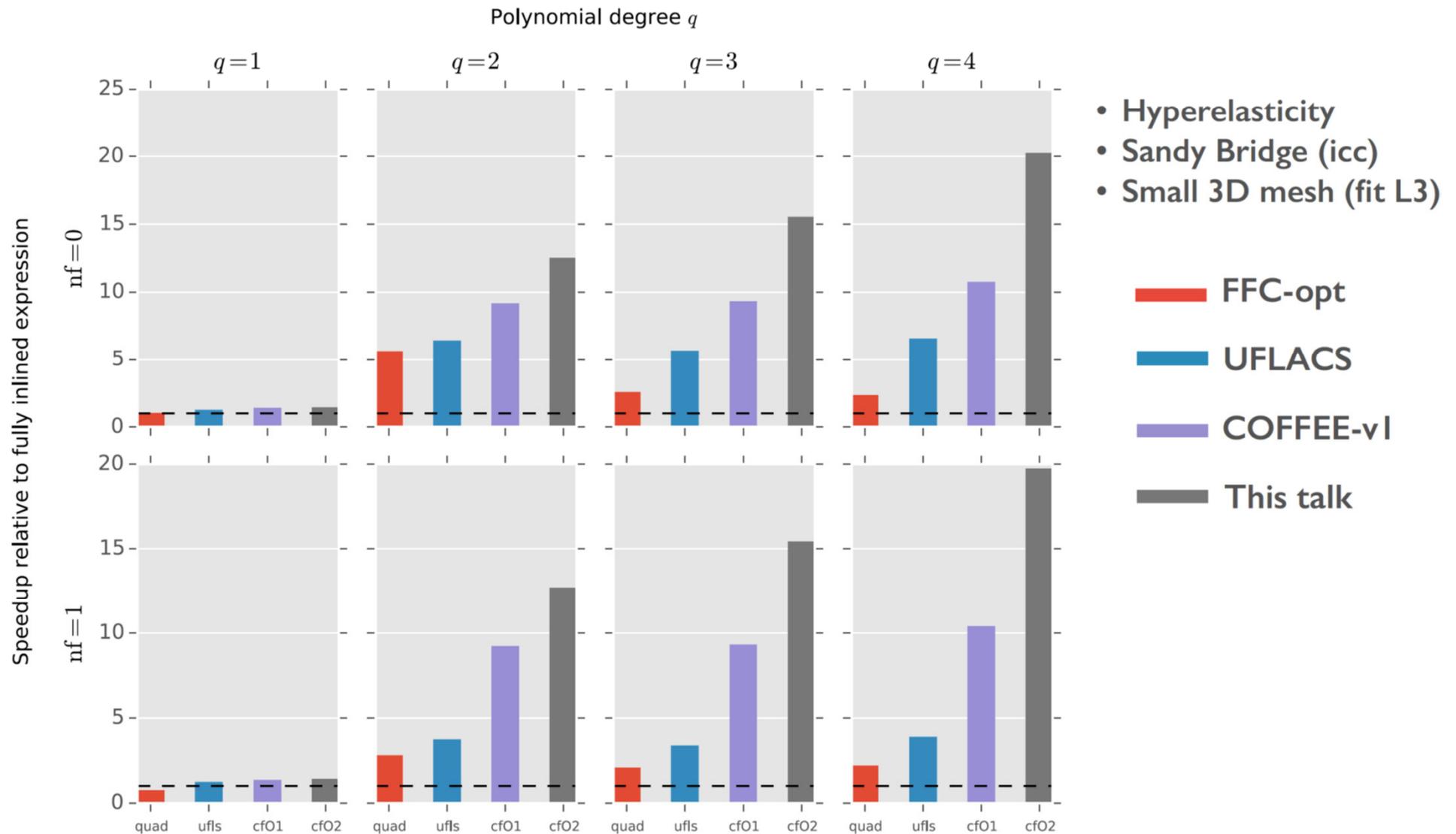
ARSENAL FOR REDUCING FLOPS

Loop-invariant code motion
Common sub-expressions elimination  **flops**



We formulate an ILP problem to find the best factorisation strategy

FOCUS ON HYPERELASTICITY



F. Luporini, D.A. Ham, P.H.J. Kelly. An algorithm for the optimization of finite element integration loops. ACM Transactions on Mathematical Software (TOMS), 2017).