

Why hybrid systems work so well for simulation based sciences

Thomas C. Schulthess





BECAUSE OF AMDAHL'S LAW!

We knew this all along, so this is trivial

But how exactly does this show in applications?

We have 20 minutes to look at two applications (that are taken from real life simulations)

Electronic structure problem (DFT-based) in material science

Collaborators: Stan Tomov¹, Azam Haidar¹, Raffaele Solcà², Antron Koszevnikov², Jack Dongarra¹

Regional model for climate simulations and weather forecasting

Collaborators: Oliver Fuhrer³, Tobias Gysi⁴, Will Sawyer⁵, David Müller⁴, Uli Schätler⁶, Michael Baldauf⁶,
Xavier Lapillone², Ugo Vareto⁵, Mauro Bianco⁵, Isabelle Bey², Tim Schröder⁷, Tom Bradley⁷

Many insightful discussion with Steve Scot⁷, John Levesque⁸, and Alessandro Curioni⁹ (and many others)

(1) ICL/UTK; (2) ETH; (3) MeteoCH; (4) SCS; (5) CSCS; (6) DWD; (7) NVIDIA; (8) Cray; (9) IBM-ZRL

Solving Kohn-Sham equation is the bottleneck of most DFT based materials science codes

Kohn-Sham Eqn.
$$\left(-\frac{\hbar^2}{2m} \nabla^2 + v_{\text{LDA}}(\vec{r}) \right) \psi_i(\vec{r}) = \epsilon_i \psi_i(\vec{r})$$

Ansatz
$$\psi_i(\vec{r}) = \sum_{\mu} c_{i\mu} \phi_{\mu}(\vec{r})$$

Hermitian matrix
$$H_{\mu\nu} = \int \phi_{\mu}^*(\vec{r}) \left(-\frac{\hbar^2}{2m} \nabla^2 + v_{\text{LDA}}(\vec{r}) \right) \phi_{\nu}(\vec{r}) d\vec{r}$$

Basis is not orthogonal
$$S_{\mu\nu} = \int \phi_{\mu}^*(\vec{r}) \phi_{\nu}(\vec{r}) d\vec{r}$$

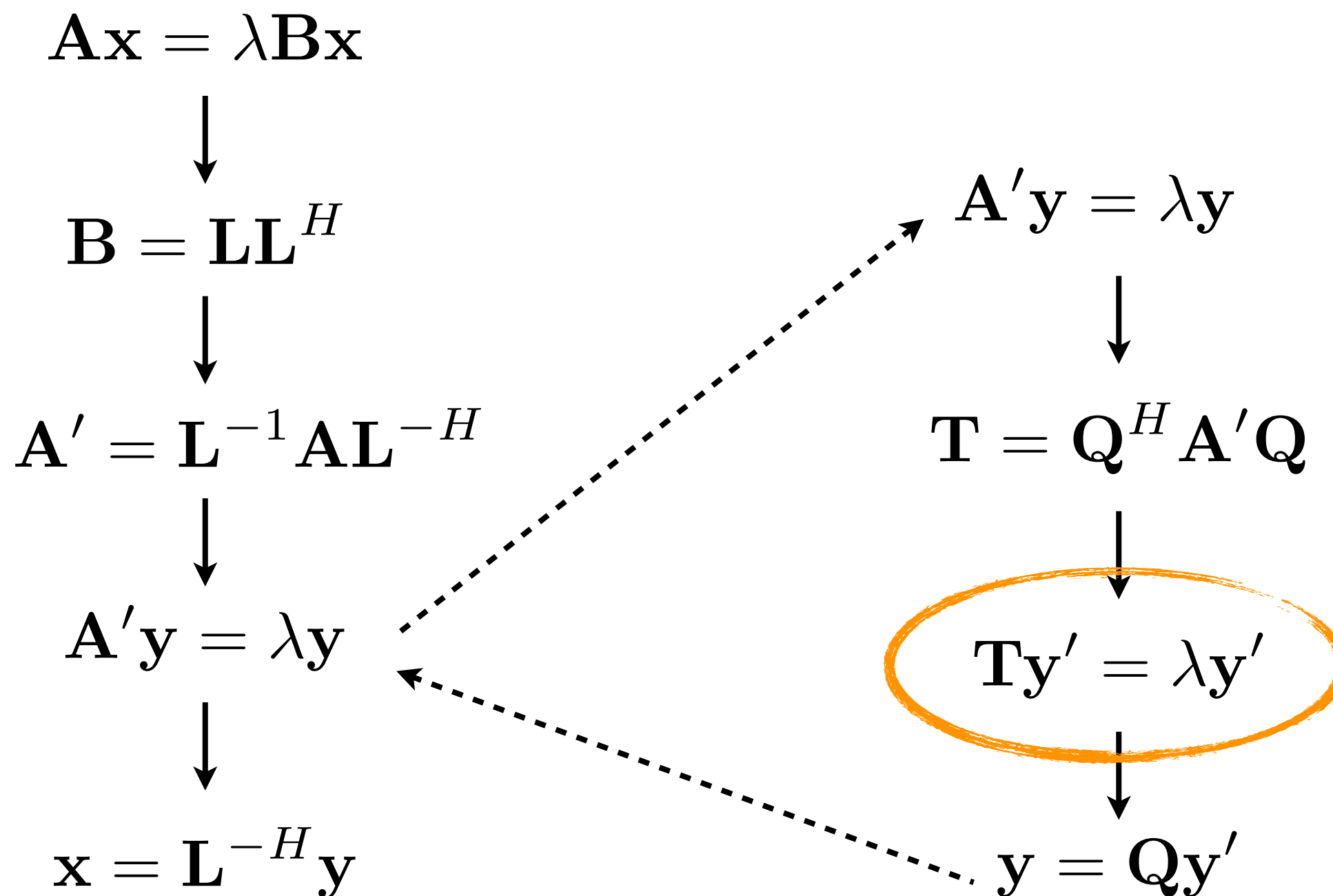
Solve generalized eigenvalue problem
$$(\mathbf{H} - \epsilon_i \mathbf{S}) = 0$$

where we are usually interested in about 10-50% of spectrum

We need eigenvectors as well, to compute the density:

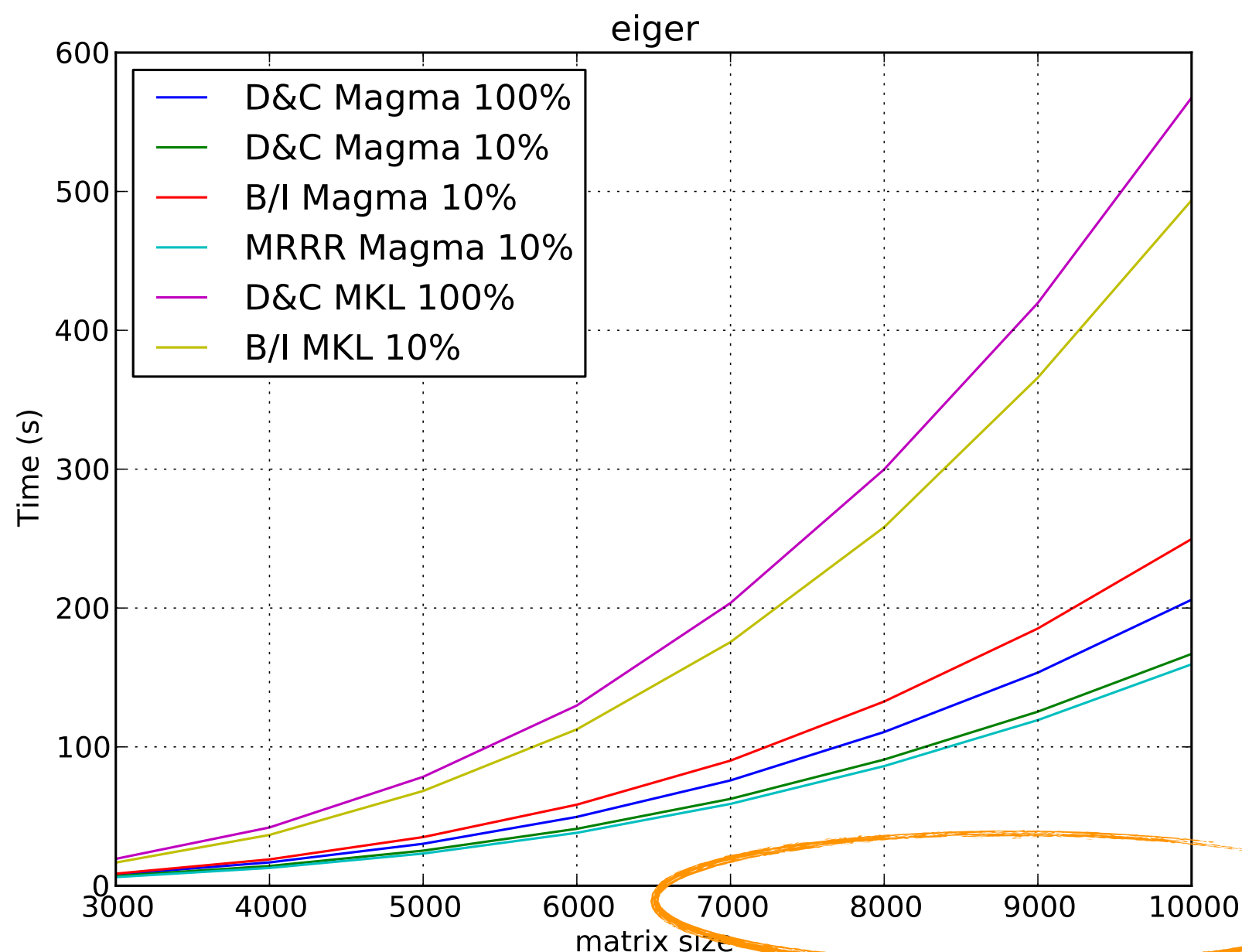
$$n(\xi) = \sum_{i=1}^N \psi_i^*(\xi) \psi_i(\xi)$$

Solving the generalized eigenvalue problem



Comparing multi-core vs. hybrid system

Test system: two 6-core AMD Opteron 2427@2.2GHz and one Tesla C2070

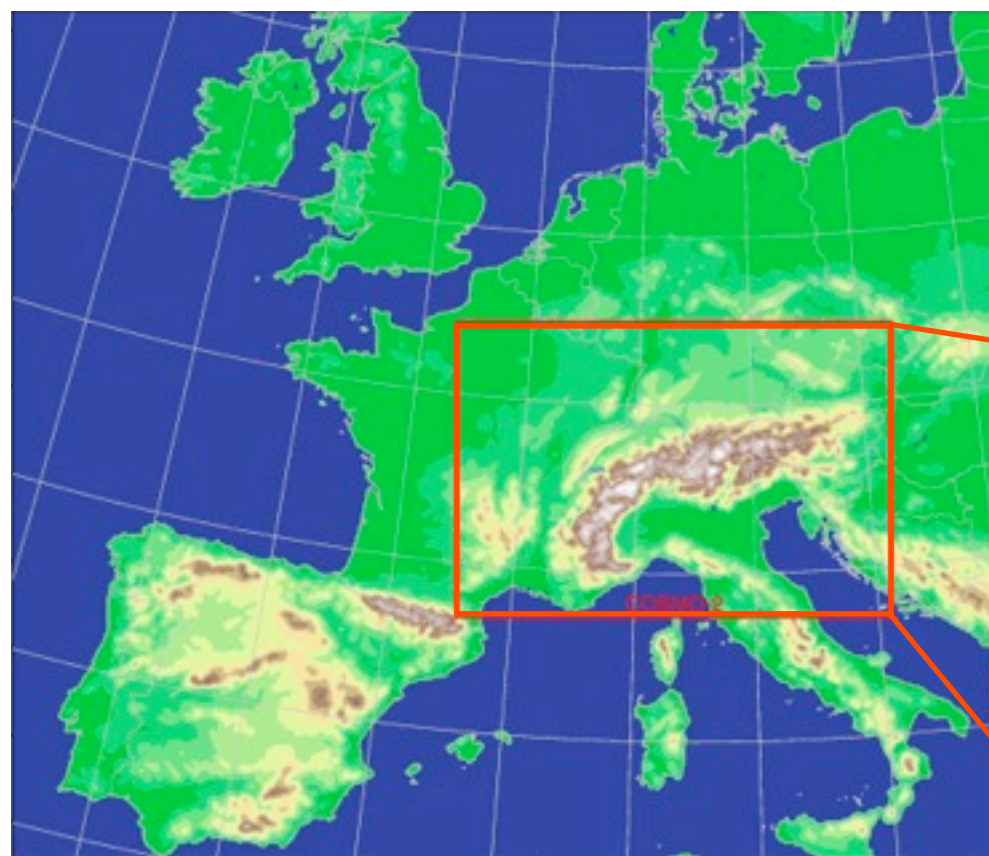


MKL runs on two AMD socket (i.e. 12 CPU cores) **MAGMA** runs on 1 CPU core and the GPU

(much care was taken to run both cases optimally)

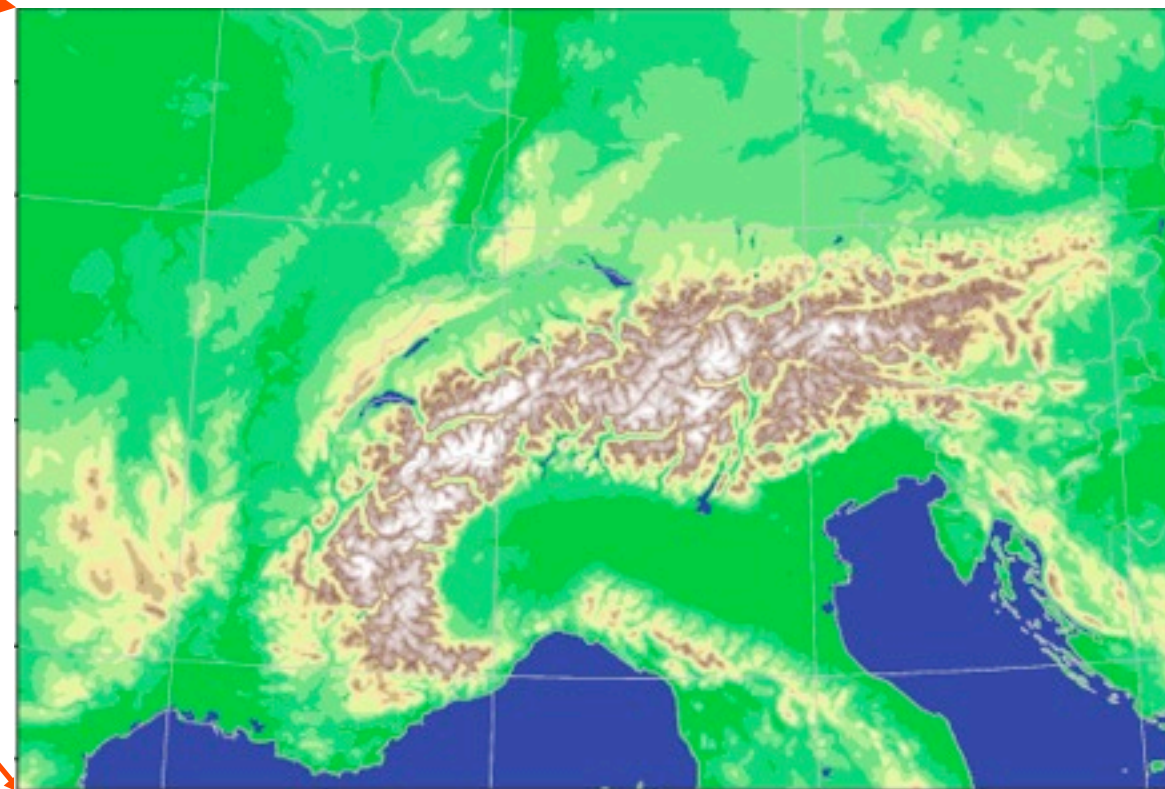
Numerical weather predictions for Switzerland (MeteoSwiss)

COSMO-7 : 3x per day 72h forecast
6.6 km lateral grid, 60 layers



ECMWF: boundary conditions
16km, 91 layers
2x per day

COSMO-2: 8x per day 24h forecast
2.2 km lateral grid , 60 layer

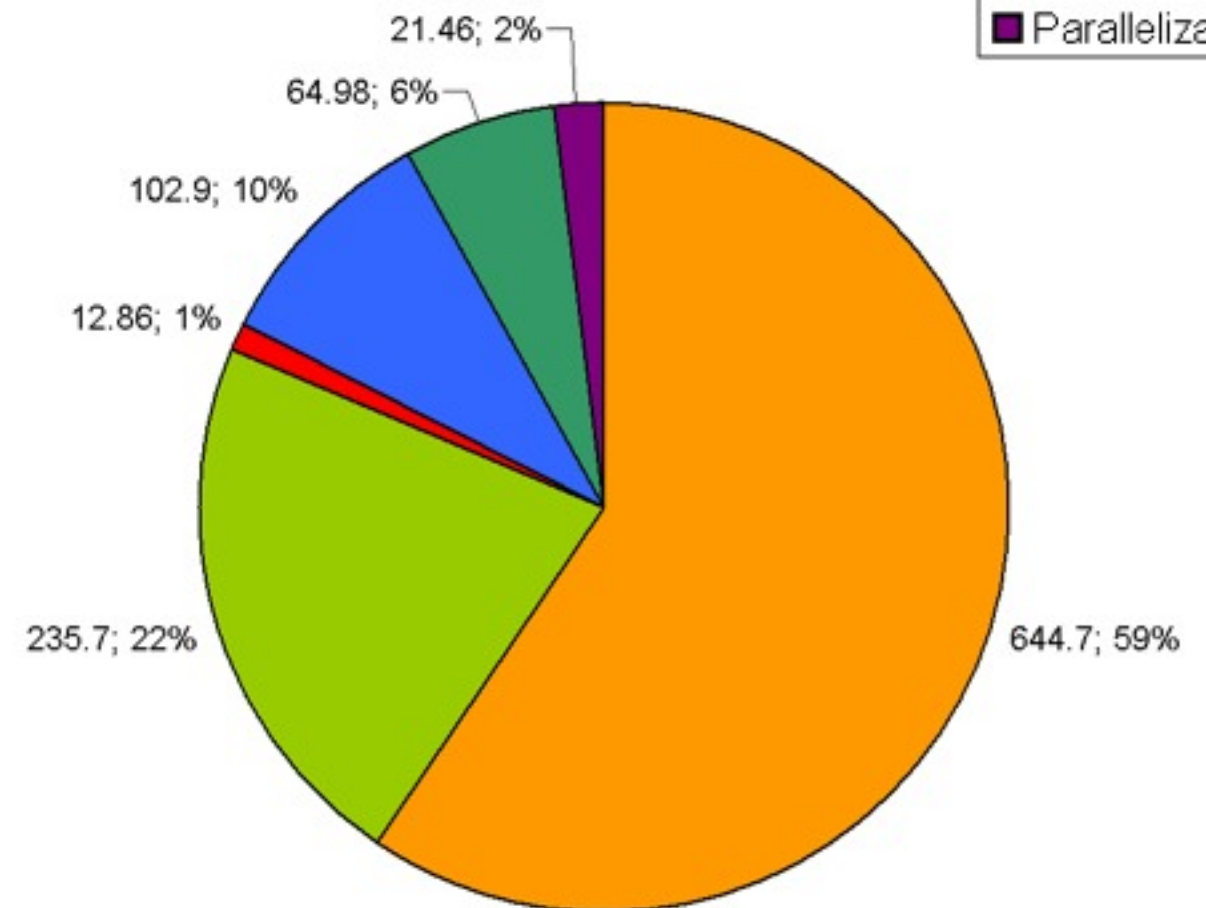
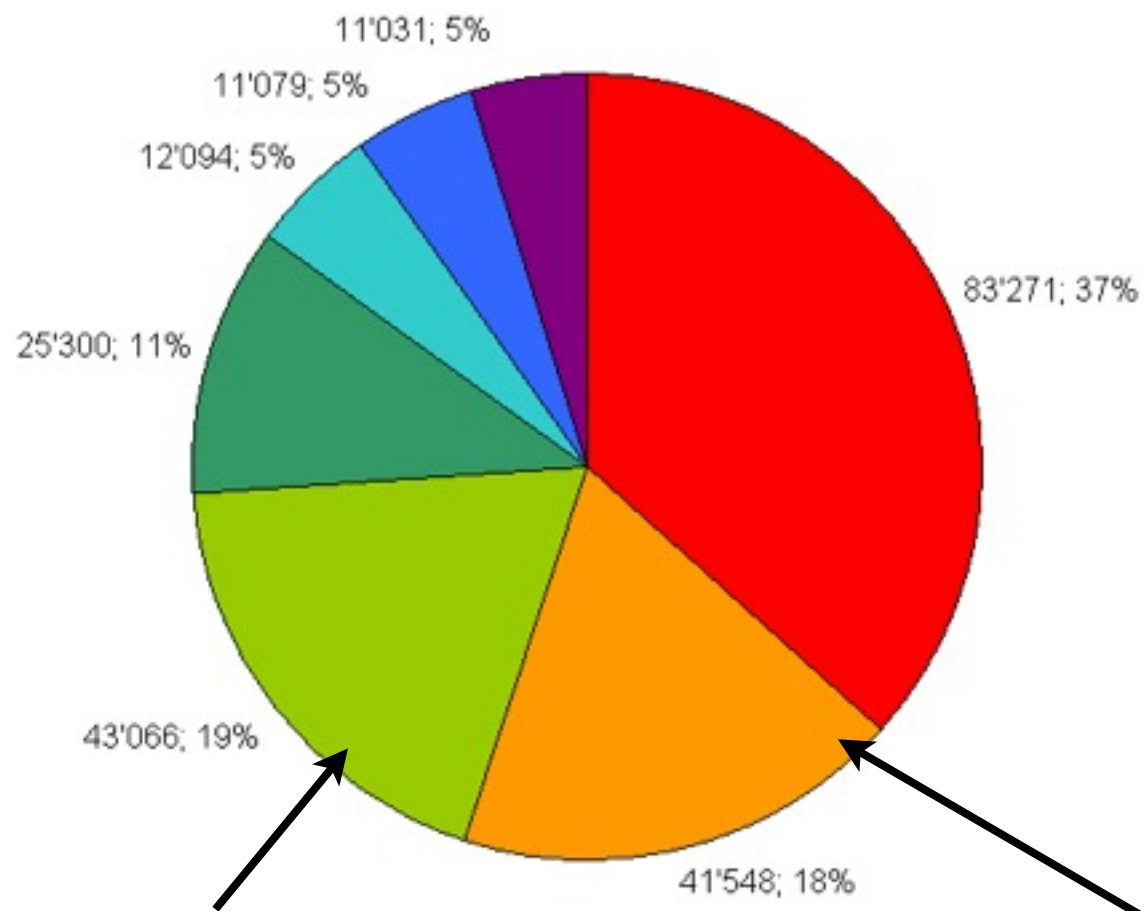


Performance profile of (original) COSMO-CCLM

Runtime based 2 km production model of MeteoSwiss

% Code Lines (F90)

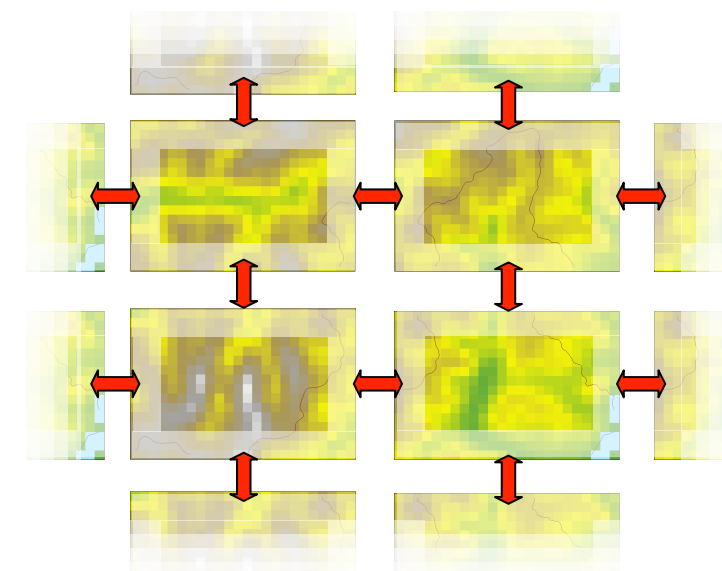
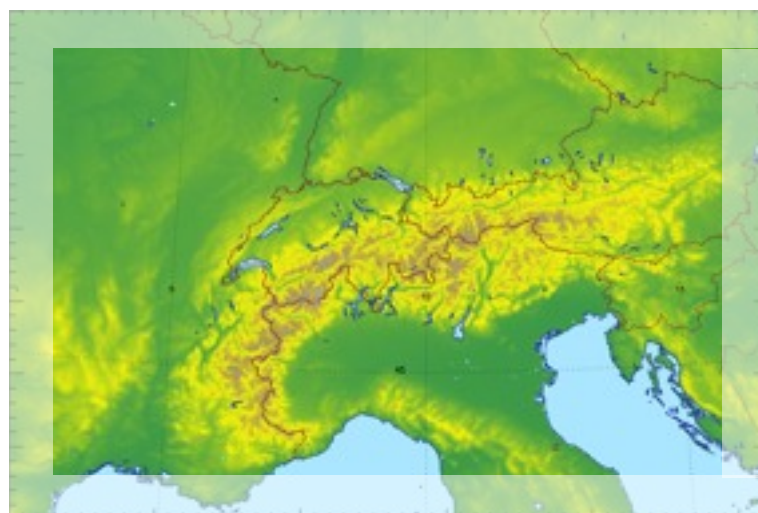
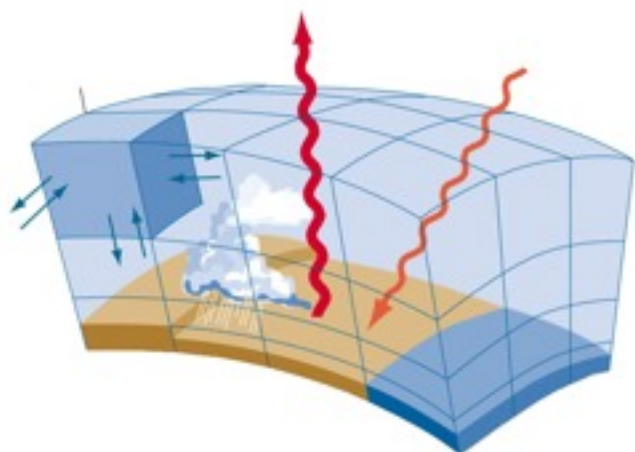
% Runtime



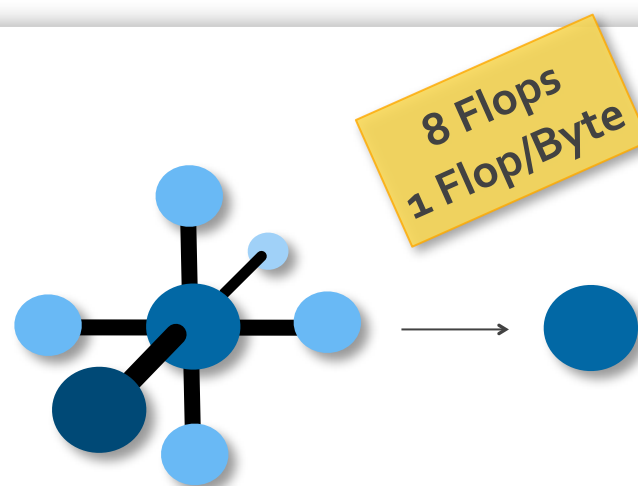
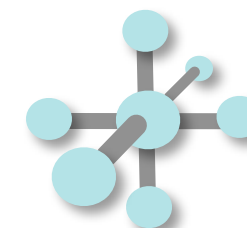
Original code (with OpenACC)

Rewrite in C++ (with CUDA)

Finite difference stencil to solve PDE on structured grid – climate, meteorology, seismic imaging, combustion

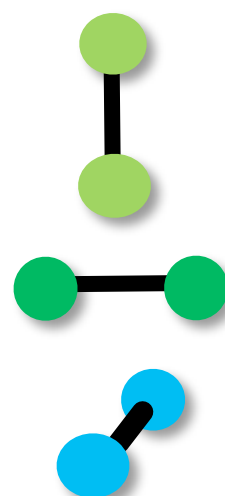


Partial differential equations on structured grid $\frac{\partial u}{\partial t} - \alpha \nabla^2 u = 0$



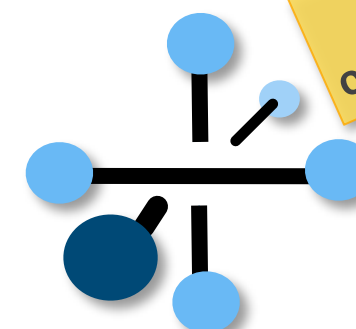
Laplacian

8 Flops
1 Flop/Byte



Divergence

8 Flops
0.5 Flops/Byte



Gradient

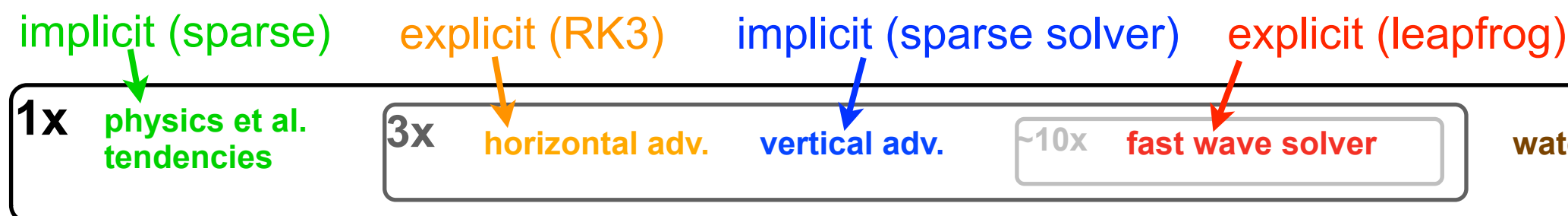
6 Flops
0.38 Flops/Byte

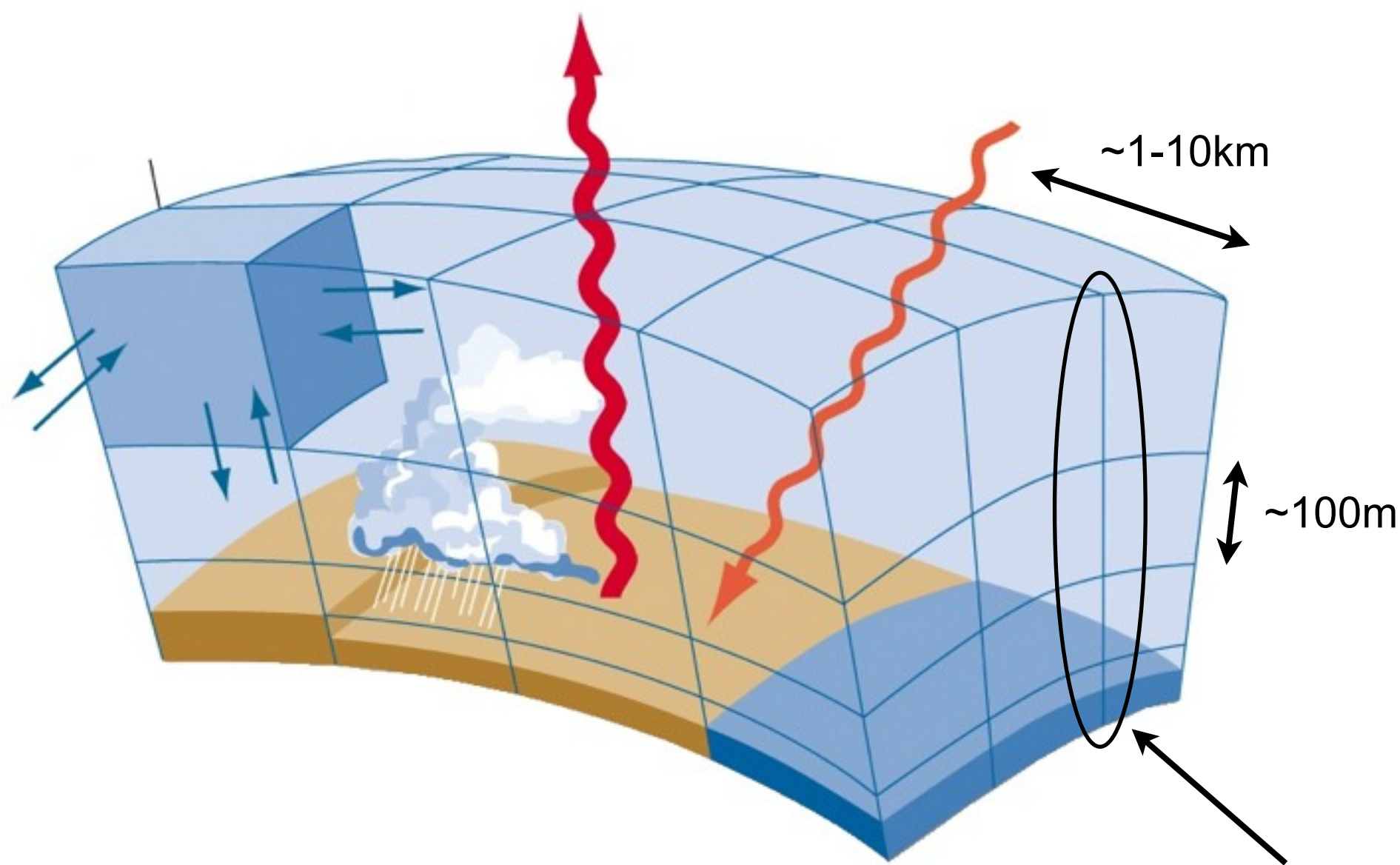
Dynamics in COSMO-CCLM

velocities	$\frac{\partial u}{\partial t} = - \left\{ \frac{1}{a \cos \varphi} \frac{\partial E_h}{\partial \lambda} - v V_a \right\} - \zeta \frac{\partial u}{\partial \zeta} - \frac{1}{\rho a \cos \varphi} \left(\frac{\partial p'}{\partial \lambda} - \frac{1}{\sqrt{\gamma}} \frac{\partial p_0}{\partial \lambda} \frac{\partial p'}{\partial \zeta} \right) + M_u$
	$\frac{\partial v}{\partial t} = - \left\{ \frac{1}{a} \frac{\partial E_h}{\partial \varphi} + u V_a \right\} - \zeta \frac{\partial v}{\partial \zeta} - \frac{1}{\rho a} \left(\frac{\partial p'}{\partial \varphi} - \frac{1}{\sqrt{\gamma}} \frac{\partial p_0}{\partial \varphi} \frac{\partial p'}{\partial \zeta} \right) + M_v$
	$\frac{\partial w}{\partial t} = - \left\{ \frac{1}{a \cos \varphi} \left(u \frac{\partial w}{\partial \lambda} + v \cos \varphi \frac{\partial w}{\partial \varphi} \right) \right\} - \zeta \frac{\partial w}{\partial \zeta} + \frac{g}{\sqrt{\gamma}} \frac{\rho_0}{\rho} \frac{\partial p'}{\partial \zeta} + M_w + g \frac{\rho_0}{\rho} \left\{ \frac{(T - T_0)}{T} - \frac{T_0 p'}{T p_0} + \left(\frac{R_v}{R_d} - 1 \right) q^v - q^l - q^f \right\}$
pressure	$\frac{\partial p'}{\partial t} = - \left\{ \frac{1}{a \cos \varphi} \left(u \frac{\partial p'}{\partial \lambda} + v \cos \varphi \frac{\partial p'}{\partial \varphi} \right) \right\} - \zeta \frac{\partial p'}{\partial \zeta} + g \rho_0 w - \frac{c_{pd}}{c_{vd}} p D$
temperature	$\frac{\partial T}{\partial t} = - \left\{ \frac{1}{a \cos \varphi} \left(u \frac{\partial T}{\partial \lambda} + v \cos \varphi \frac{\partial T}{\partial \varphi} \right) \right\} - \zeta \frac{\partial T}{\partial \zeta} - \frac{1}{\rho c_{vd}} p D + Q_T$
water	$\frac{\partial q^v}{\partial t}$
	$\frac{\partial q^{l,f}}{\partial t}$
turbulence	$\frac{\partial e_t}{\partial t} = - \left\{ \frac{1}{a \cos \varphi} \left(u \frac{\partial e_t}{\partial \lambda} + v \cos \varphi \frac{\partial e_t}{\partial \varphi} \right) \right\} - \zeta \frac{\partial e_t}{\partial \zeta} + K_m^v \frac{g \rho_0}{\sqrt{\gamma}} \left\{ \left(\frac{\partial u}{\partial \zeta} \right)^2 + \left(\frac{\partial v}{\partial \zeta} \right)^2 \right\} + \frac{g}{\rho \theta_v} F^{\theta_v} - \frac{\sqrt{2} e_t^{3/2}}{\alpha_M l} + M_{e_t}$

- > Structured grid / stencil computation is ideal for GPU (as a bandwidth engine)
- > Small memory footprint allows up to put entire problem on “accelerator(s)”
- > But the tridiagonal solve is inherently a serial

Timestep





Tridiagonal solve for matrix $\sim 50-100$

Current status of implementation:

- > structure grid motif works well on data parallel, high-memory bandwidth “accelerator”
- > tridiagonal solve makes good use of cash on CPU and single thread performance

Hybrid system we need for simulation based science has to have four “things”

- Many data parallel “high throughput” processing units
- Something that gives the highest single (several) thread performance
- Enough memory to store vectors/matrices that supports highest possible bandwidth
- Lots of memory to store everything else and supports decent bandwidth

Note: I have not talked about distributed memory systems because the possibilities are vast and I only had 20 minutes



QUESTIONS / COMMENTS