

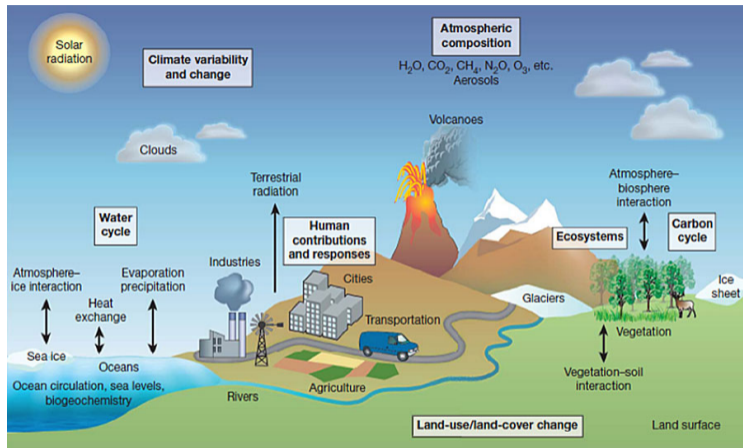
To improve accuracy of weather and climate models via a reduction of precision

Peter Düben

University Research Fellow of the Royal Society

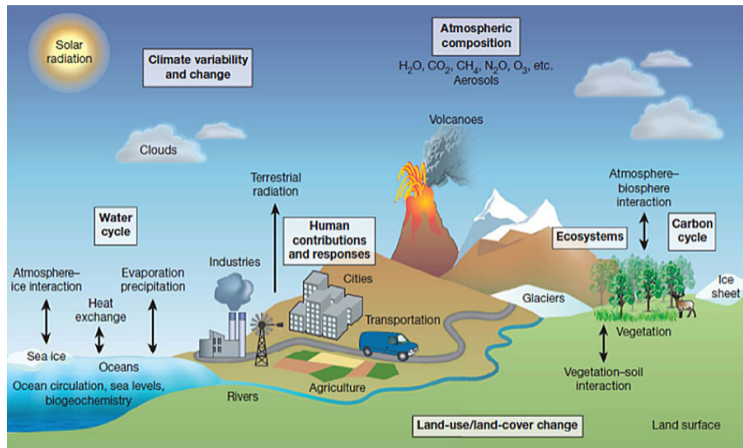
European Centre for Medium-Range Weather Forecasts (ECMWF)

Predicting weather and climate: Why is it so hard?



www.gfdl.noaa.gov

Predicting weather and climate: Why is it so hard?

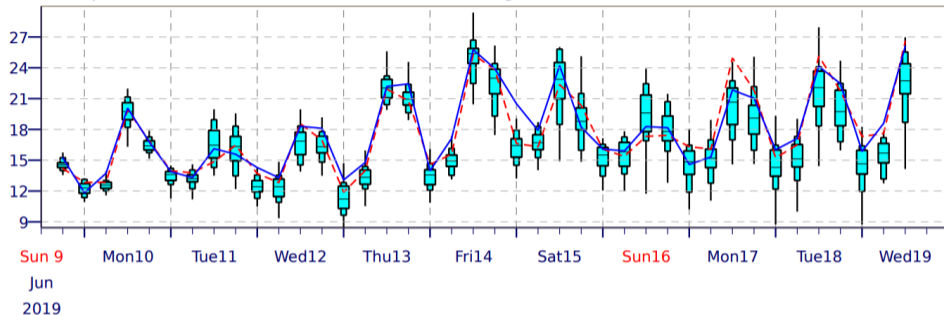


www.gfdl.noaa.gov

The Earth System is complex, huge and chaotic and we do not have sufficient resolution to resolve all important processes.

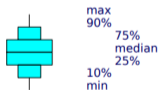
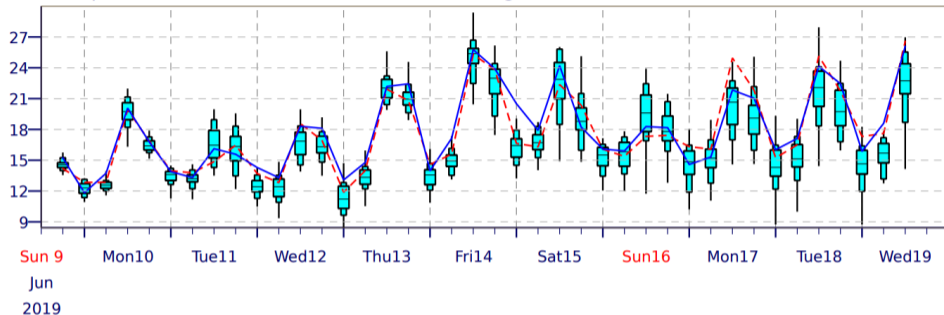
We use ensemble forecasts to prediction model uncertainty

2m Temperature(°C) reduced to 472 m (station height) from 477 m (HRES) and 382 m (ENS)



We use ensemble forecasts to prediction model uncertainty

2m Temperature(°C) reduced to 472 m (station height) from 477 m (HRES) and 382 m (ENS)



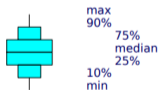
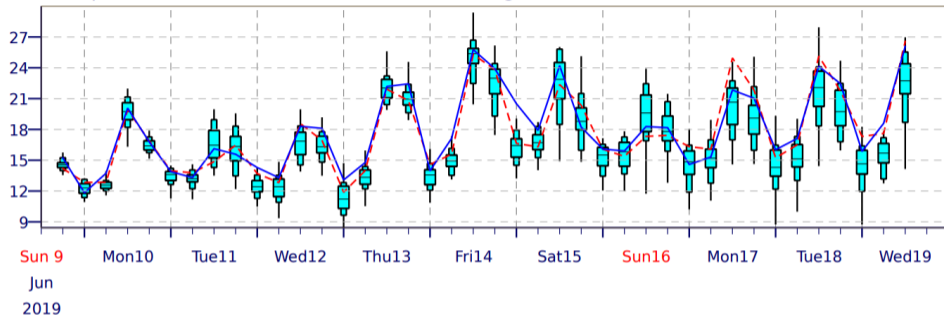
ENS Control(16 km)

High Resolution (8 km)

The magnitude of errors is predicted via the spread of ensemble forecasts.

We use ensemble forecasts to prediction model uncertainty

2m Temperature(°C) reduced to 472 m (station height) from 477 m (HRES) and 382 m (ENS)



ENS Control(16 km)

High Resolution (8 km)

The magnitude of errors is predicted via the spread of ensemble forecasts.

This spread can help us to adjust precision; error due to precision « spread is good.

Less numerical precision → more computing power

Double precision (64 bits) is used almost exclusively in weather and climate modelling.

Less numerical precision → more computing power

Double precision (64 bits) is used almost exclusively in weather and climate modelling.

Reduce numerical precision

→ lower power, higher performance.

→ higher resolution or increased complexity.

→ more accurate predictions of future weather and climate.

Less numerical precision → more computing power

Double precision (64 bits) is used almost exclusively in weather and climate modelling.

Reduce numerical precision

→ lower power, higher performance.

→ higher resolution or increased complexity.

→ more accurate predictions of future weather and climate.

Temperature in Mainz:

double precision (64 bits): 14.561192512512207°C

single precision (32 bits): 14.5611925°C

half precision (16 bits): 14.5625°C

Less numerical precision → more computing power

Double precision (64 bits) is used almost exclusively in weather and climate modelling.

Reduce numerical precision

→ lower power, higher performance.

→ higher resolution or increased complexity.

→ more accurate predictions of future weather and climate.

Temperature in Mainz:

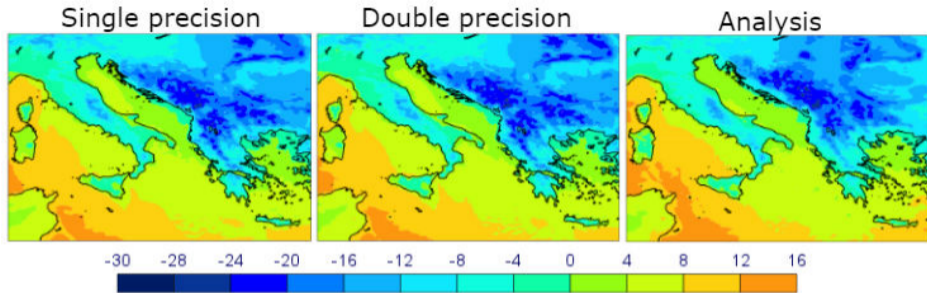
double precision (64 bits): 14.561192512512207°C

single precision (32 bits): 14.5611925°C

half precision (16 bits): 14.5625°C

But can we really do it? And how far can we go?

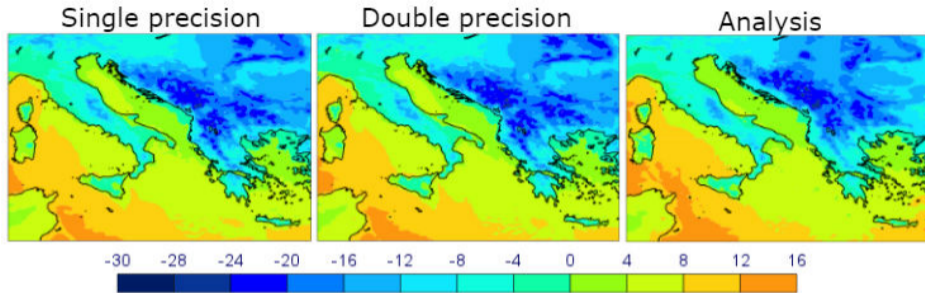
ECMWF's weather forecast model in single precision



- ▶ Forecast quality in double and single precision is almost identical.
- ▶ 40% reduction of run time.
- ▶ Benefit for global simulations at cloud-resolving resolution.

Düben and Palmer MWR 2014; Váňa, Düben et al. MWR 2017; Düben et al. ECMWF Newsletter 2018

ECMWF's weather forecast model in single precision



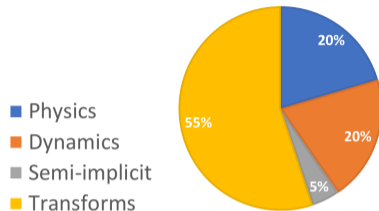
- ▶ Forecast quality in double and single precision is almost identical.
- ▶ 40% reduction of run time.
- ▶ Benefit for global simulations at cloud-resolving resolution.

Düben and Palmer MWR 2014; Váňa, Düben et al. MWR 2017; Düben et al. ECMWF Newsletter 2018

Can we go lower than single precision?

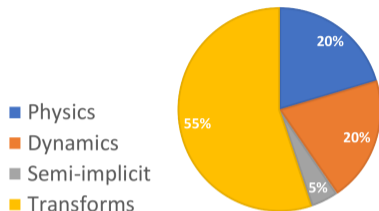
Machine learning hardware for fast simulations with low precision

Relative cost for model components for a non-hydrostatic model at 1.45 km resolution:



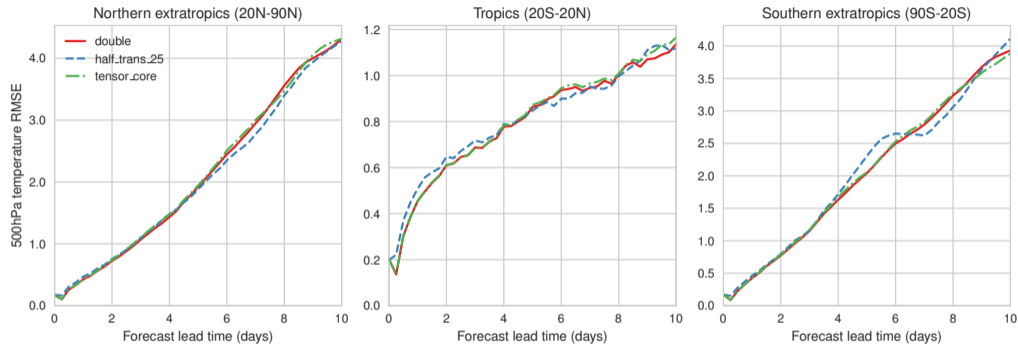
Machine learning hardware for fast simulations with low precision

Relative cost for model components for a non-hydrostatic model at 1.45 km resolution:



- ▶ The Legendre transform is the most expensive kernel. It consists of a large number of standard matrix-matrix multiplications.
- ▶ If we can re-scale the input and output fields, we can use half precision arithmetic.
- ▶ Tensor Cores on NVIDIA Volta GPUs are optimised for half-precision matrix-matrix calculations with single precision output. 7.8 TFlops for double precision vs. 125 TFlops for half precision on the Tensor Core.

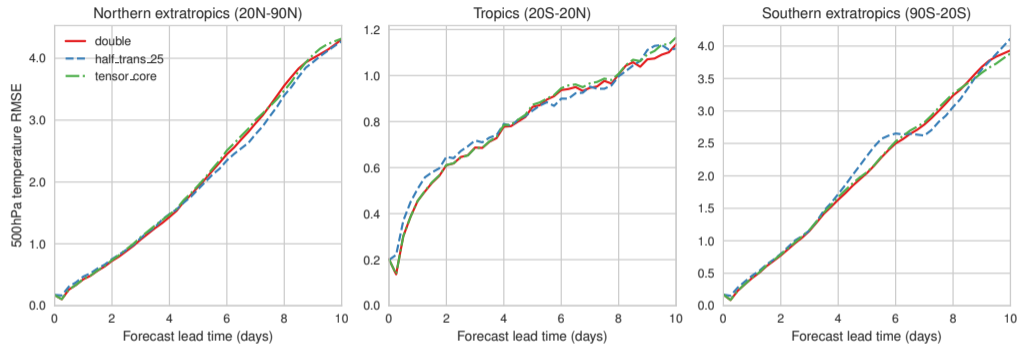
Half precision Legendre Transformations



Root-mean-square error for Z500 at 9 km resolution averaged over multiple start dates.

Hatfield, Chantry, Dueben, Palmer **Best Paper Award PASC2019**.

Half precision Legendre Transformations



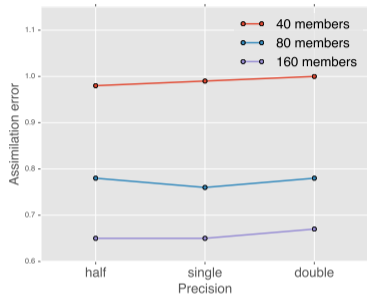
Root-mean-square error for Z500 at 9 km resolution averaged over multiple start dates.

Hatfield, Chantry, Dueben, Palmer **Best Paper Award PASC2019**.

The simulations are using an emulator to reduce precision.

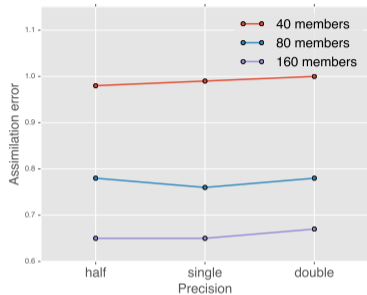
Dawson and Dueben GMD 2017

Data assimilation with reduced precision



Data assimilation in Lorenz'95 using an Ensemble Kalman filter. Hatfield, Dueben, Palmer JAMES 2018

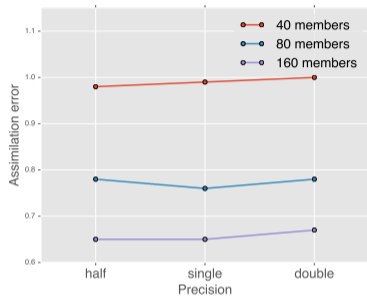
Data assimilation with reduced precision



Data assimilation in Lorenz'95 using an Ensemble Kalman filter. Hatfield, Dueben, Palmer JAMES 2018

A large ensemble at low precision is better than a small ensemble at high precision.

Data assimilation with reduced precision

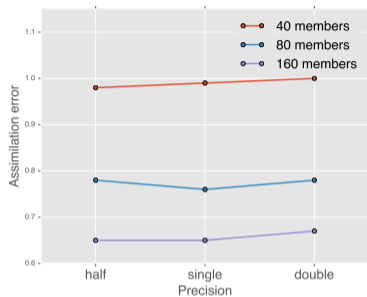


Data assimilation in Lorenz'95 using an Ensemble Kalman filter. Hatfield, Dueben, Palmer JAMES 2018

A large ensemble at low precision is better than a small ensemble at high precision.

We gain almost one “day” in terms of predictability.

Data assimilation with reduced precision



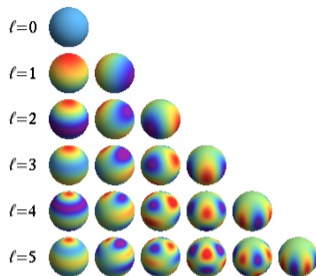
Data assimilation in Lorenz'95 using an Ensemble Kalman filter. Hatfield, Dueben, Palmer JAMES 2018

A large ensemble at low precision is better than a small ensemble at high precision.

We gain almost one "day" in terms of predictability.

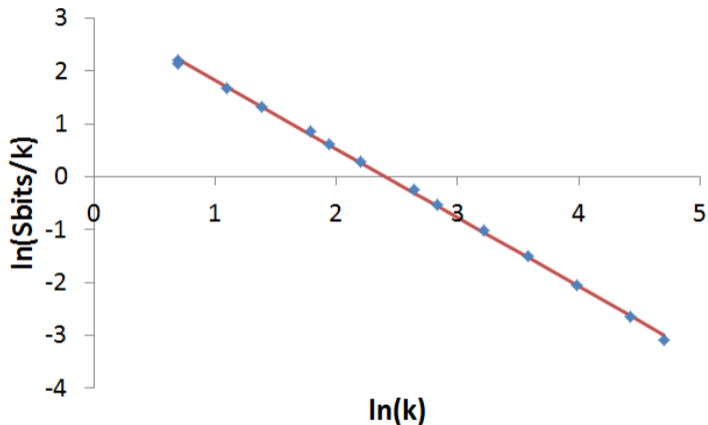
However, 4DVar data assimilation may be more difficult...

A scale-selective approach



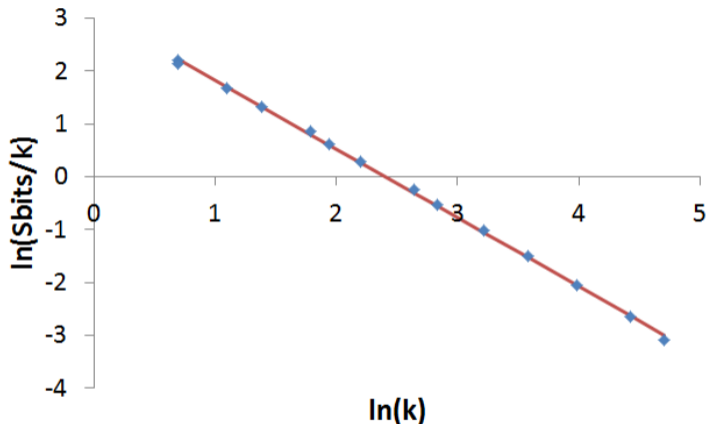
- ▶ Spectral models allow to treat different scales at different precision.
- ▶ We can reduce precision when calculating the small scales.
- ▶ This is intuitive due to the high inherent uncertainty in small scale dynamics (parametrisation, viscosity, data-assimilation,...).
- ▶ The smallest scales are most expensive.

A scale-selective approach



A scale-dependent reduction in precision for the surface quasi-geostrophic equations.

A scale-selective approach

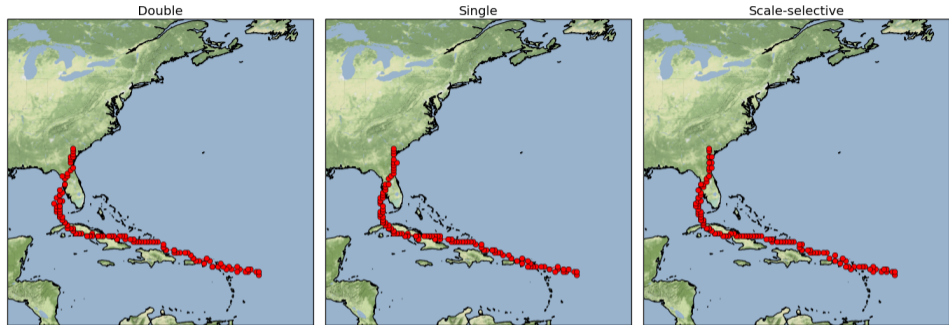


A scale-dependent reduction in precision for the surface quasi-geostrophic equations.

Forecast simulations confirm that a scale-selective approach is much more efficient than a uniform precision reduction.

Thornes, Düben and Palmer QJRMS 2017, Thornes, Düben and Palmer QJRMS 2018

A scale-selective approach: Track of Hurricane Irma



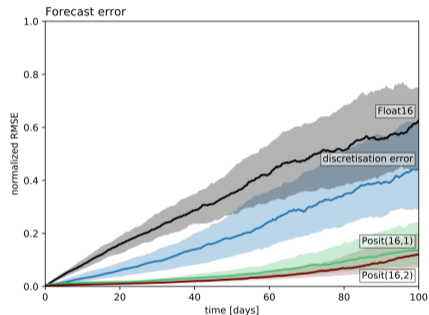
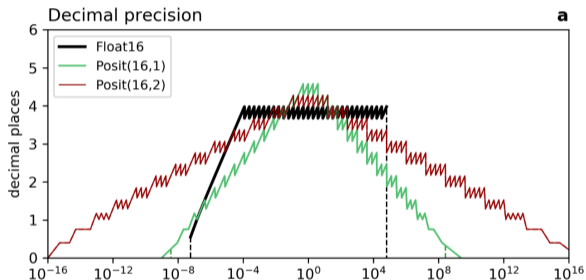
- ▶ Simulations with OpenIFS at 40 km resolution.
- ▶ The scale-selective simulation is using scale-selective precision in spectral space. An average of 8.6 bits is used for the significant.

What number format to use?

16 bits is not much so you may need to show some flexibility and use Posits.

What number format to use?

16 bits is not much so you may need to show some flexibility and use Posits.

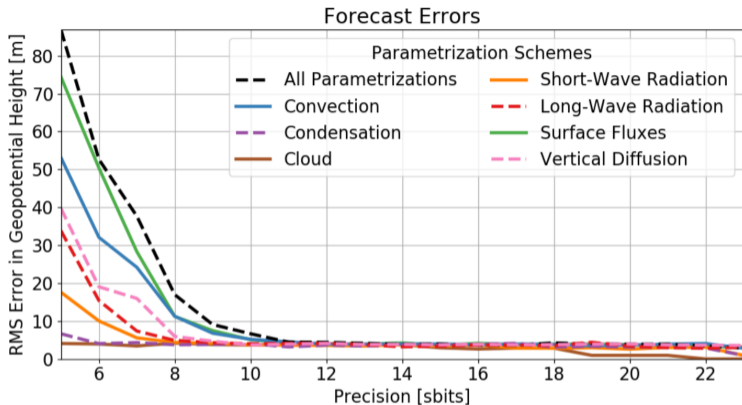


Left: Dynamic number ranges of 16 bit Posit formats and 16 bit half precision floats.

Right: Forecast error for a shallow water model if reduced precision is used.

Kloewer, Düben and Palmer CONGA 2019

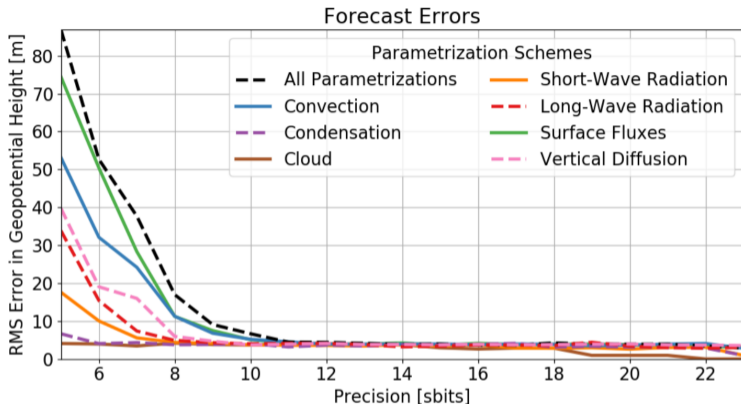
Reduced precision in physical parametrisation schemes



Saffin, Hatfield, Düben and Palmer in prep.

Forecast error with respect to double precision at 2-weeks lead time if numerical precision is reduced in the different parametrisation schemes in the SPEEDY model.

Reduced precision in physical parametrisation schemes

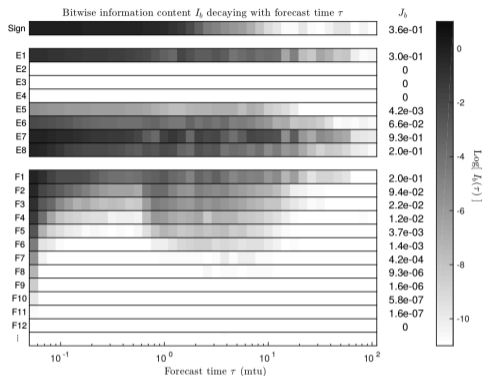


Saffin, Hatfield, Düben and Palmer in prep.

Forecast error with respect to double precision at 2-weeks lead time if numerical precision is reduced in the different parametrisation schemes in the SPEEDY model.

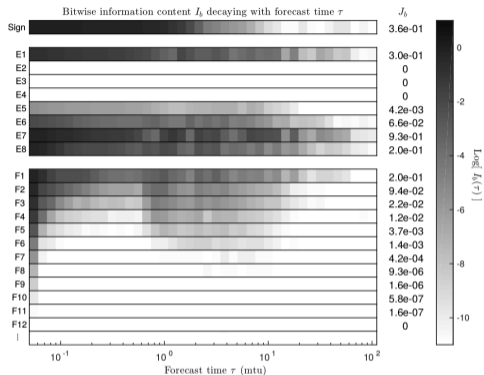
There is scope!

Bitwise information content and predictability



Information content of bits for a Lorenz'63 model using a single long term integration and Shannon information theory.

Bitwise information content and predictability



Information content of bits for a Lorenz'63 model using a single long term integration and Shannon information theory.

It is possible to identify information content of individual bits and their impact on predictability into the future.

Reduce precision in weather and climate models

What we already know:

- ▶ We can reduce precision significantly in many model components (forecast model - Dueben and Palmer 2014; data assimilation - Hatfield et al. 2018; land surface - Dawson et al. 2017; ocean models - Tinto et al. 2019).
- ▶ We can emulate reduced numerical precision in large Fortran models (Dawson and Dueben 2017).
- ▶ We can adjust numerical precision to forecast accuracy and model uncertainty (Dueben et al. 2018).
- ▶ We can provide a physical justification for numerical precision from domain knowledge (e.g. via scale selective precision - Chantry et al. 2018).
- ▶ We can use Shannon Information Theory to extract the information content per bit (Jeffress et al. 2017).

Reduce precision in weather and climate models

What we still need:

- ▶ Tools that allow an automated search for the optimal precision level when non-linear feedbacks are present.
- ▶ A basic understanding how to formulate models to minimize numerical precision (re-scaling of equations, perturbation approaches, multi-grid solvers...).
- ▶ Tools to predict a performance increase from a precision reduction for a given hardware.
- ▶ Information how future hardware and hardware co-design will look like (CPUs, GPUs, TPUs, FPGAs, ASICs...).

Conclusions

- ▶ Reducing precision can free resources to increase resolution of weather and climate models.
- ▶ Single precision is providing almost identical forecast skill when compared to double precision simulations.
- ▶ A detailed performance analysis is important to understand the benefit of model developments.
- ▶ For single precision, savings are mainly generated via a reduction of cash misses and improved vectorization.
- ▶ A further reduction beyond single precision for expensive kernels is possible and promising.
- ▶ We will need better performance models to drive precision reduction in the future.



Funded by the
European Union

The presenter gratefully acknowledges the funding a Royal Society University Research Fellowship and from the ESIWACE2 project. The ESIWACE2 project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 823988.