

Lowering precision in an atmospheric ensemble data assimilation system

Sam Hatfield (1), Peter Düben (2), Matthew Chantry (1), Tim Palmer (1)

samuel.hatfield@physics.ox.ac.uk, (1) University of Oxford (2) European Centre for Medium-Range Weather Forecasts



Introduction

Conventional wisdom dictates that, in atmospheric models, it is always best to use the highest numerical precision available. Only recently, several studies have found a significant tolerance in models to a reduction in precision, and therefore a potential free source of computational resources. The aim of this project is to extend these investigations into data assimilation.

1. Motivation: reducing precision to improve data assimilation

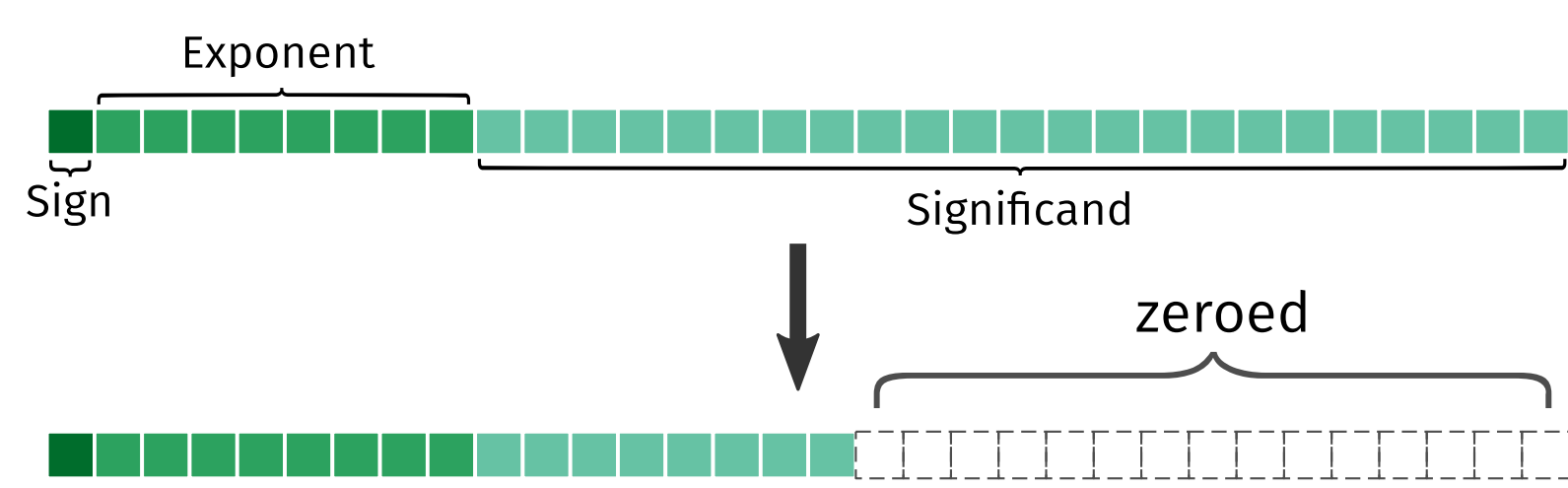


Figure 1 Procedure for emulating a reduced-precision floating-point number; in this case the trailing 15 bits of the significand are set to 0.

- Model simulations, including those used for data assimilation, are **typically carried out at double precision**: using 64 bits per variable
- However, this may be **more precision than is necessary** given our uncertainties about the model formulation and the noisiness and sparseness of the observations
- If we can reduce precision without impacting the quality of data assimilation, we can reinvest any saved resources into, for example, the ensemble size
- This would hypothetically **improve the quality of assimilation for no extra cost**
- In this study, we investigate how the quality of data assimilation is affected by a reduction in numerical precision
- We **emulate future reduced-precision hardware** with an emulator capable of fully IEEE-754 compliant half-precision

2. Methods: model and assimilation setup

- We use an intermediate complexity GCM, **SPEEDY**, to generate a nature run
- From this nature run we **derive observations** by **adding random noise** to the variables
- To assimilate observations, we use the **local ensemble transform Kalman filter (LETKF)** with SPEEDY as the assimilation model
- We use 20 members for assimilation with the **relaxation-to-prior perturbations (RTPP)** covariance inflation scheme and Gaspari-Cohn covariance localisation
- We use the **software emulator** to lower the precision of model variables **below single precision**
- We compare **64 bit SPEEDY** with **22 bit SPEEDY** (10 bit significand)
- The **22 bit model has clear biases**, with respect to the 64 bit model (Figure 3)
- However, it may be that this **reduced-precision error is within the margin of uncertainty** provided by model error and observation error
- Our objective is to measure how a substantial reduction in numerical precision in the forecast model affects the quality of the analysis

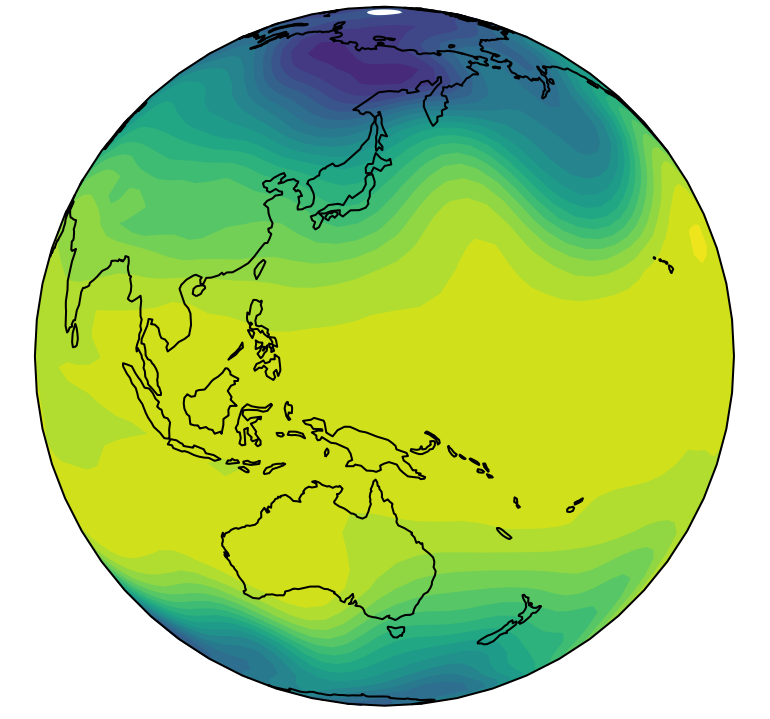


Figure 2 Snapshot of a SPEEDY simulation.

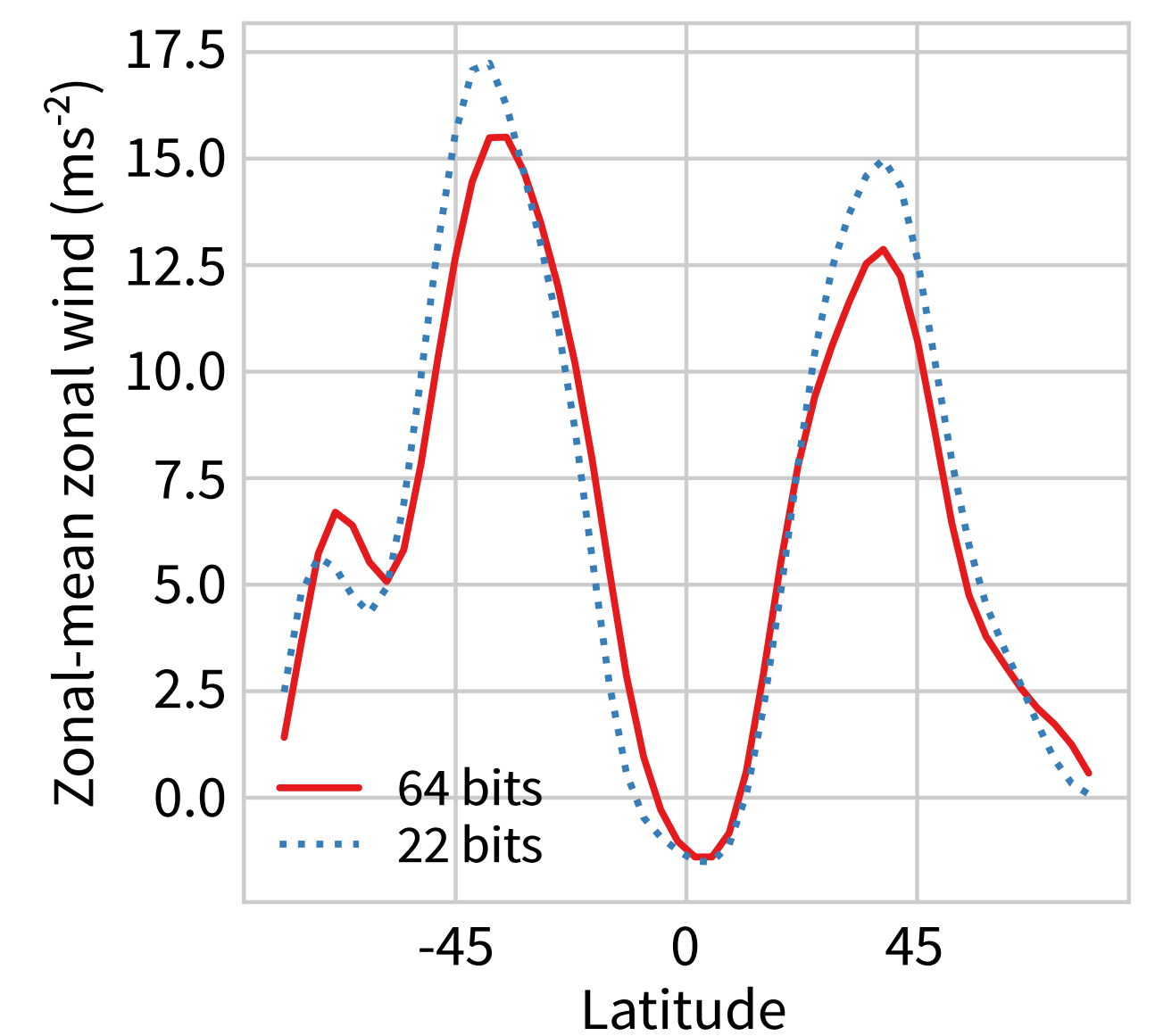


Figure 3 The zonal wind for the 64 bit and 22 bit models.

3. Results: model error vs. rounding error

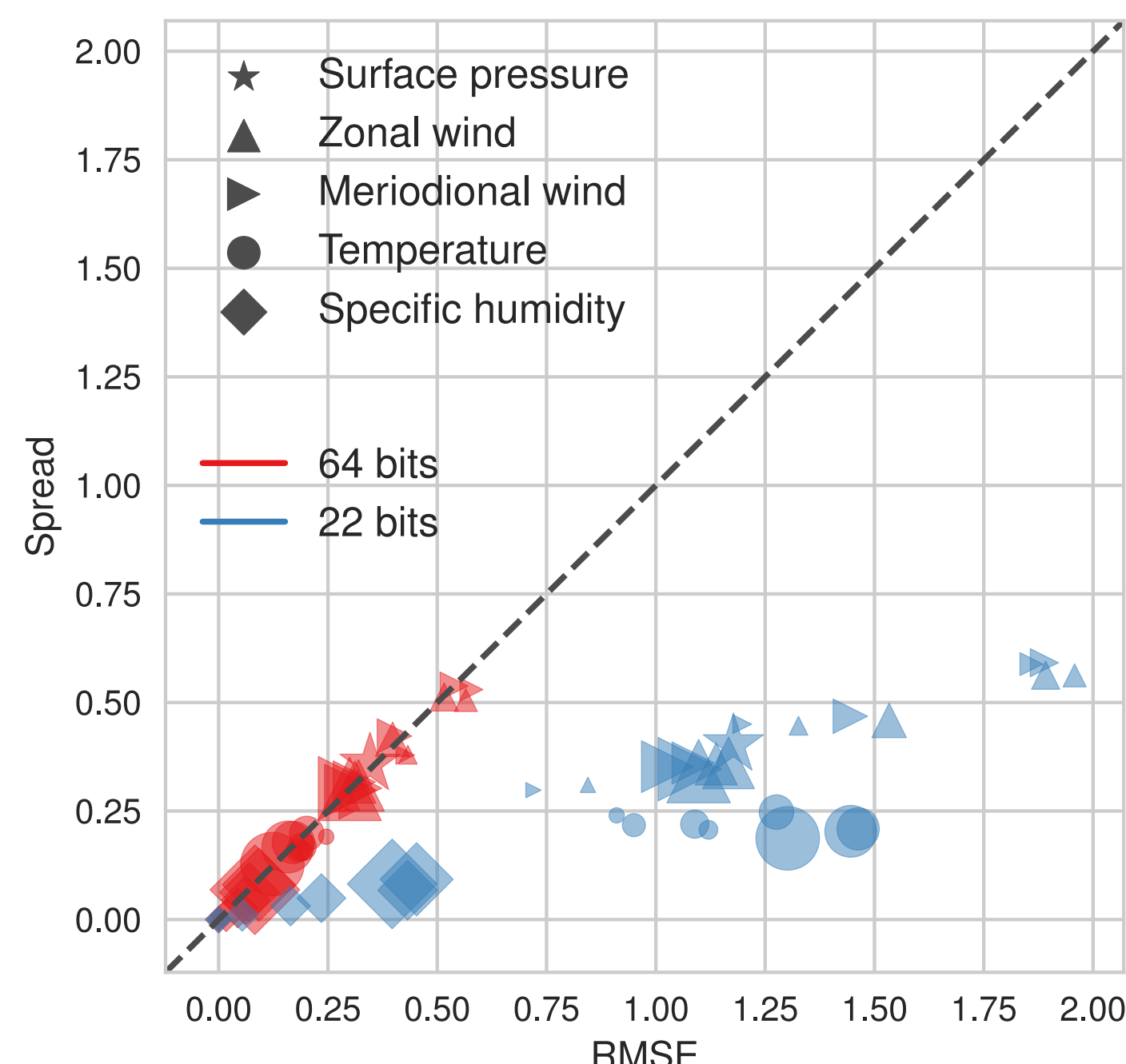


Figure 4 The error and spread for the 64 bit and 22 bit model setups in a perfect model scenario. The marker size is proportional to the model level height.

- Figure 4 compares the analysis error and spread of the 64 bit and 22 bit setups
- It looks like the 22 bit model **performs substantially worse** than the 64 bit model
- However, the nature run is also 64 bits, so the two setups cannot be fairly compared

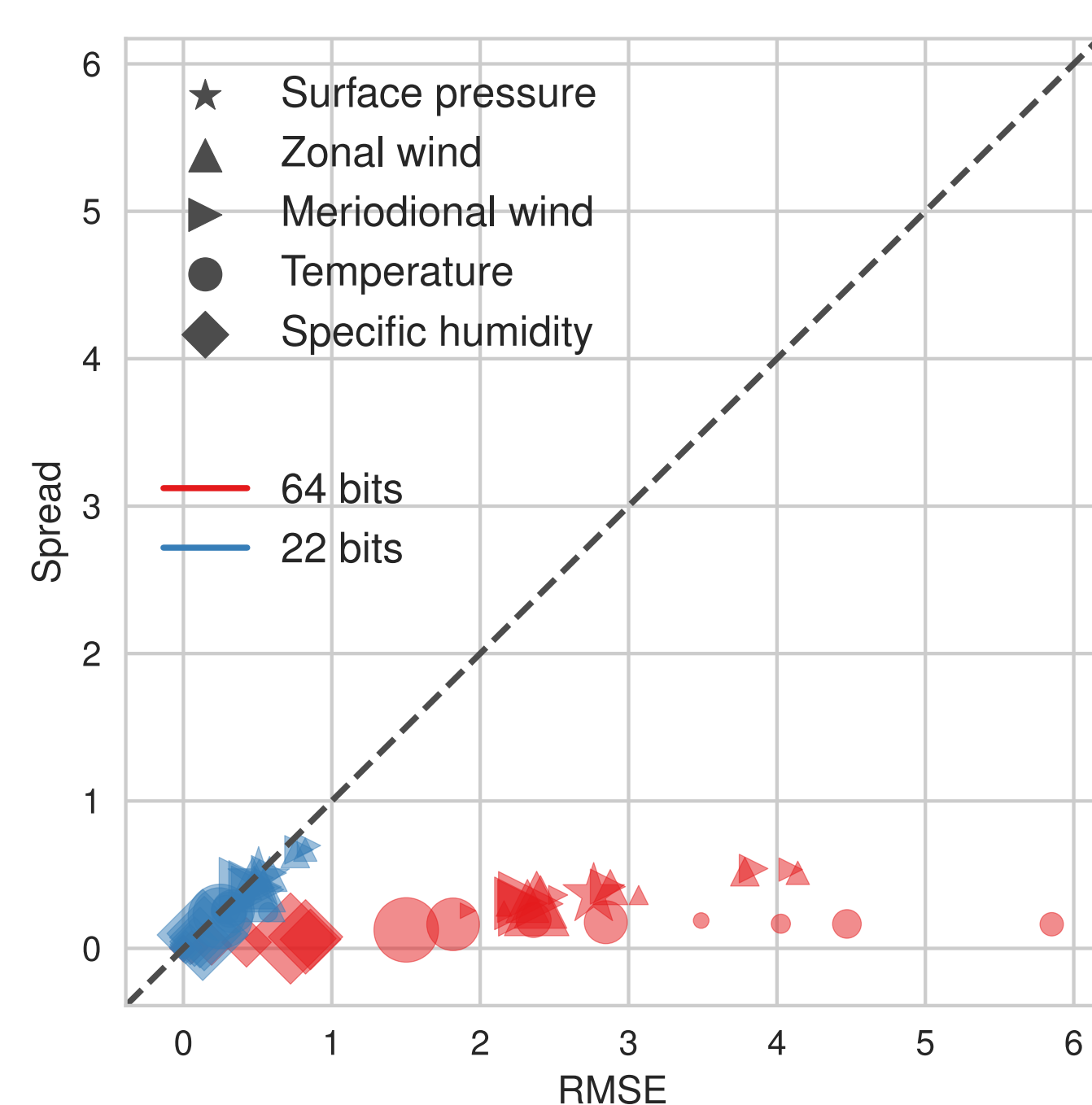


Figure 5 The same as Figure 4, but the nature run is generated using the 22 bit model.

- Figure 5 shows the same experiment as Figure 4, but the nature run was generated using the 22 bit model
- In this case, **the results are swapped**
- The apparent degradation in Figure 4 caused by the precision reduction is therefore simply a consequence of how the nature run was generated

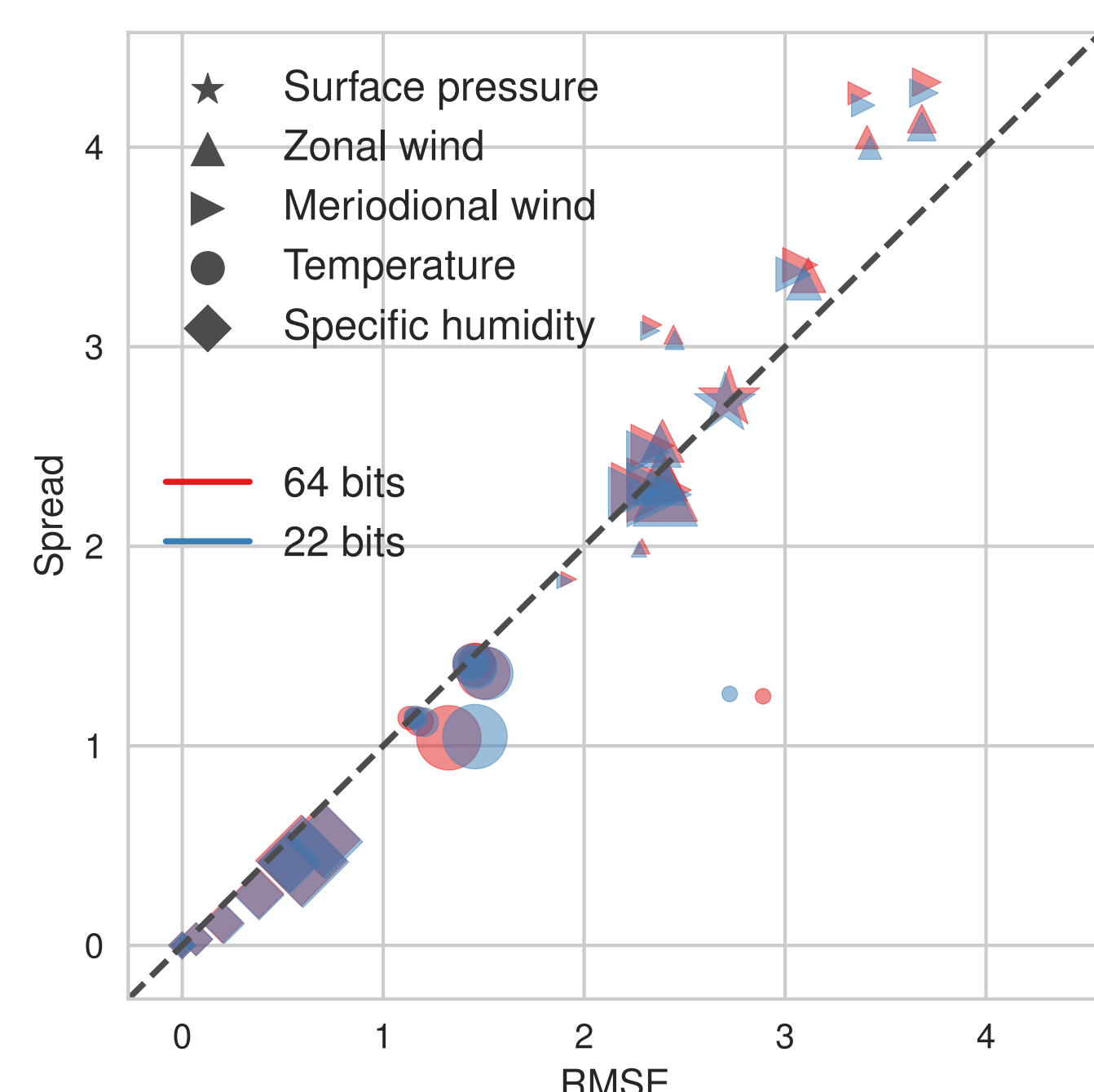


Figure 6 The same as Figure 4, but the nature run is generated using a higher resolution (T39) model.

- In order to **fairly compare the 64 bit and 22 bit models**, we need some model error
- Figure 6 shows the same experiment as Figure 4, but **the nature run was generated using a higher resolution - T39** instead of T30
- The gap in analysis quality between the 22 bit and 64 bit models is almost eliminated
- This is because the **model error due to unresolved scales "hides" the error from reducing precision**
- Even a modest amount of model error can hide large rounding errors

4. Ongoing work: half precision spectral transforms

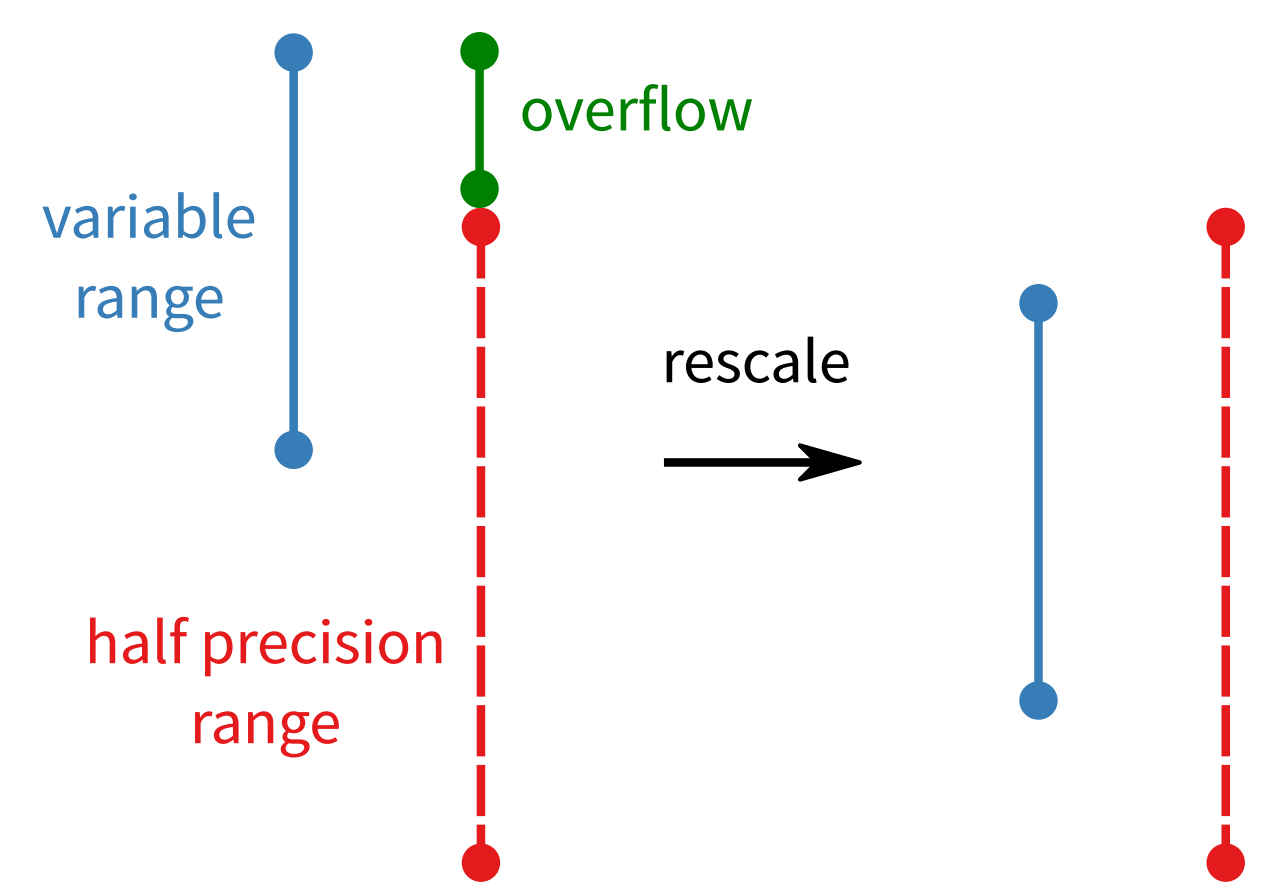


Figure 7 A simple re-scaling procedure which allows us to use half-precision computations.

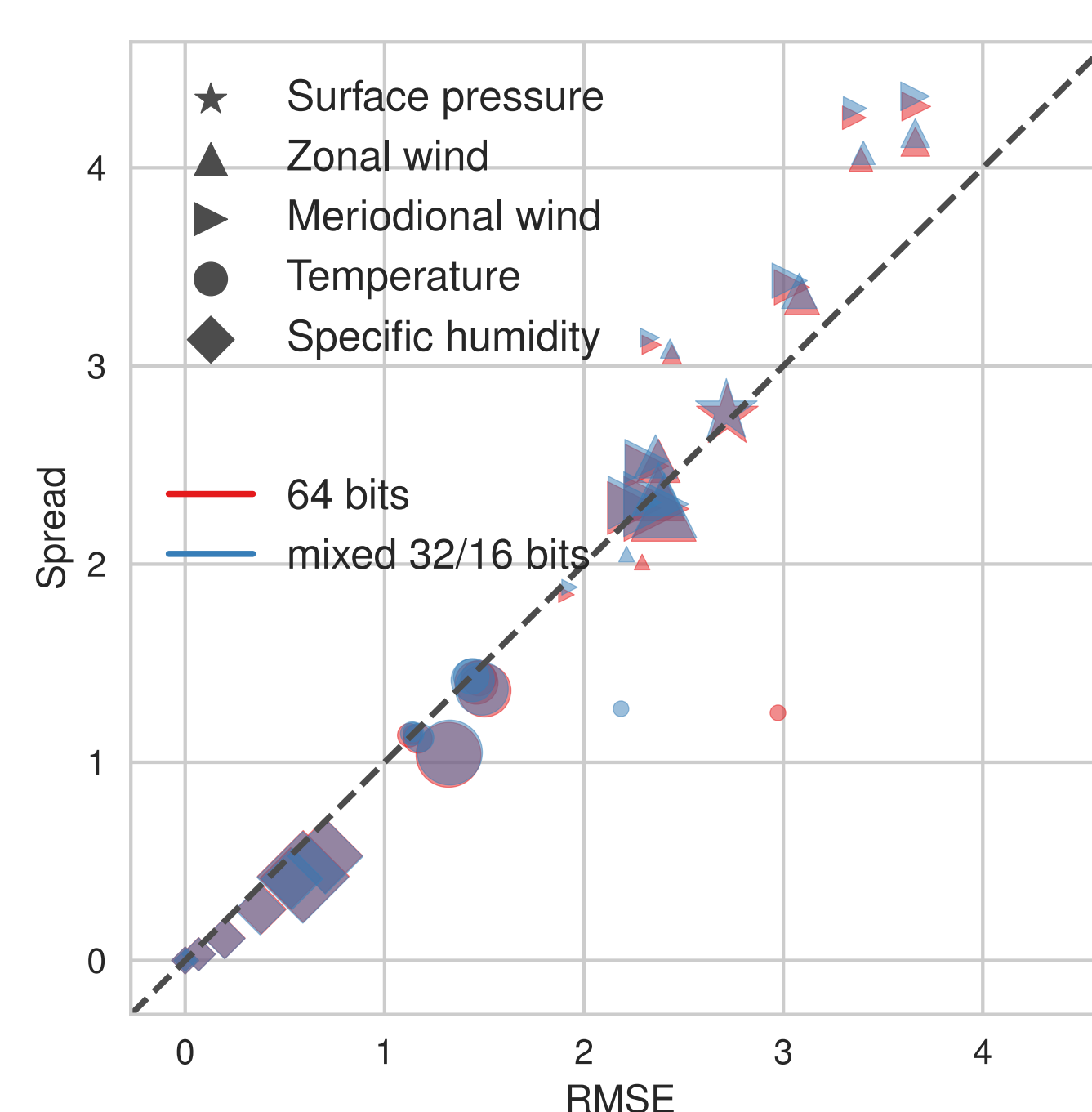


Figure 8 The same as Figure 6, but a 32/16 bit "mixed precision" model is compared with the 64 bit model.

- So far we have considered only a reduction in the significand width (see Figure 1)
- In order to consider half-precision (16 bit) computations, **we must also reduce the exponent width**, to only 5 bits
- This **reduces the range of representable numbers**, to only around 11 orders of magnitude, which is not enough for most model variables
- However, with a simple trick, we can use **half-precision for the expensive spectral transforms** (Figure 7)
- This results in a **mixed-precision model** - half precision for the spectral transforms, single precision for everything else
- This model performs comparably to the 64 bit model in an imperfect model data assimilation experiment (Figure 8)

5. Conclusion and outlook

- In the presence of model error precision can be reduced **substantially from the double precision standard with minimal impact on assimilation quality**
- Future work will focus further on using half-precision arithmetic within atmospheric models
- Modern GPUs such as the Nvidia Volta allow a **factor of 16 acceleration of half precision computations** with respect to double-precision
- Hardware with half-precision support will increasingly become available to high performance computing users in the future (see e.g. Piz Daint, Swiss National Supercomputing Centre)

Thanks to Takemasa Miyoshi and Keiichi Kondo for help with the SPEEDY experiments.

