# Exascale data archive: A new way of processing high-resolution climate data

Frédéric Laliberté[†], Paul J. Kushner

Department of Physics, University of Toronto, Canada

[†]Corresponding author: frederic.laliberte@utoronto.ca

UNIVERSITY OF TORONTO

## A. Motivation

The vast amount of climate data made available from new simulations and recent observations promises to revolutionize climate science. However, without an effective strategy to manage the $O(10\,\text{Pb})$ CMIP5/CORDEX database and the $O(1000\,\text{Pb})$ of its successor, progress could be stalled. In the future, data storage costs will be dominated by energy costs so that in 10 years a $O(1000\,\text{Pb})$ archive will be 3 times more expensive (constant dollars) to maintain than the actual $O(10\,\text{Pb})$ CMIP5/CORDEX.[1] This suggests a future where data redundancies will be unaffordable and where users will only store processed data on their local storage array.

## Advanced Climate Diagnostics

In order to allow the development of better diagnostics, it is imperative that any data management strategy includes provisions for data-intensive analyses. In higher resolution simulations, the use of increasingly complex diagnostics is often necessary to quantify teleconnections and multiscale processes. Here, we present an example of a recent diagnostic with expansive data requirements.
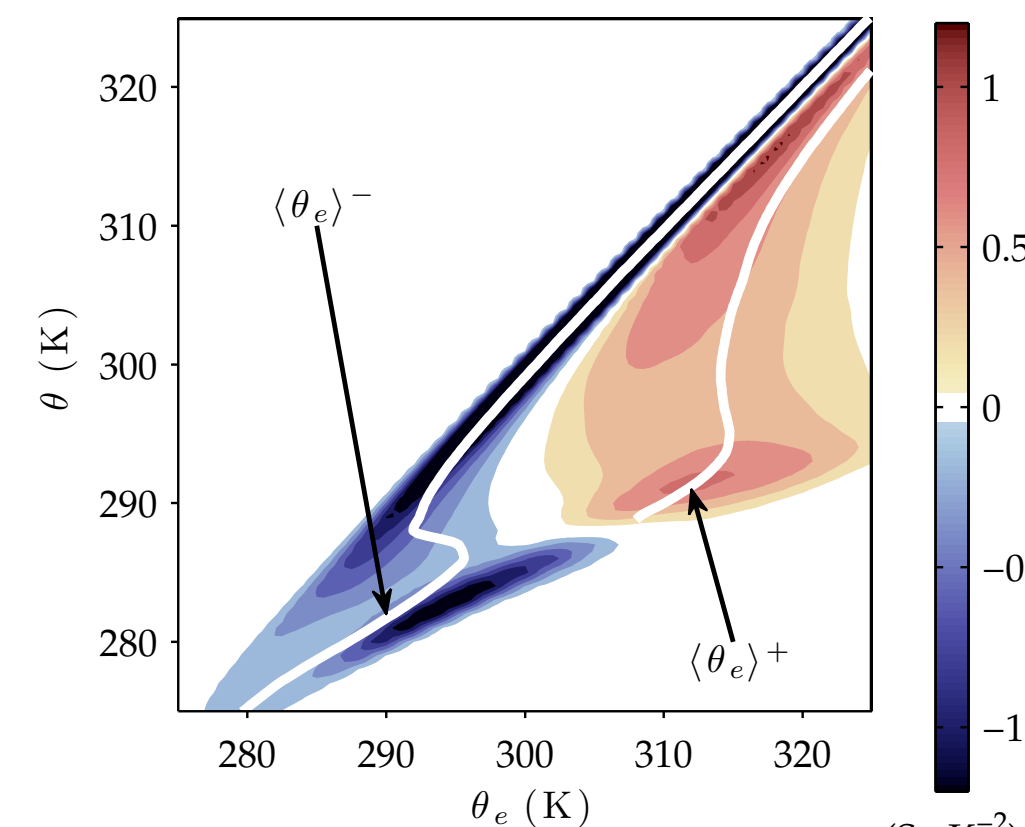
## Example: Mass Flux Joint Distribution

The Mass Flux Joint Distribution[2], $M$, quantifies meridional mass fluxes according to both their potential temperature $\theta$ and equivalent potential temperatures $\theta_e$. These two quantities are different when moisture is present in which case they are related by:

$$\theta_e \approx \theta e^{\frac{L_v}{c_p T} q_T}.$$

The joint distribution is given by:

$$M(\phi, \theta', \theta'_e) = \frac{a\cos\phi}{2\pi T}\int_0^T \int_0^{2\pi}\int_0^{P_s} v\delta(\theta-\theta')\delta(\theta_e-\theta'_e)\frac{dp}{g}d\lambda dt.$$



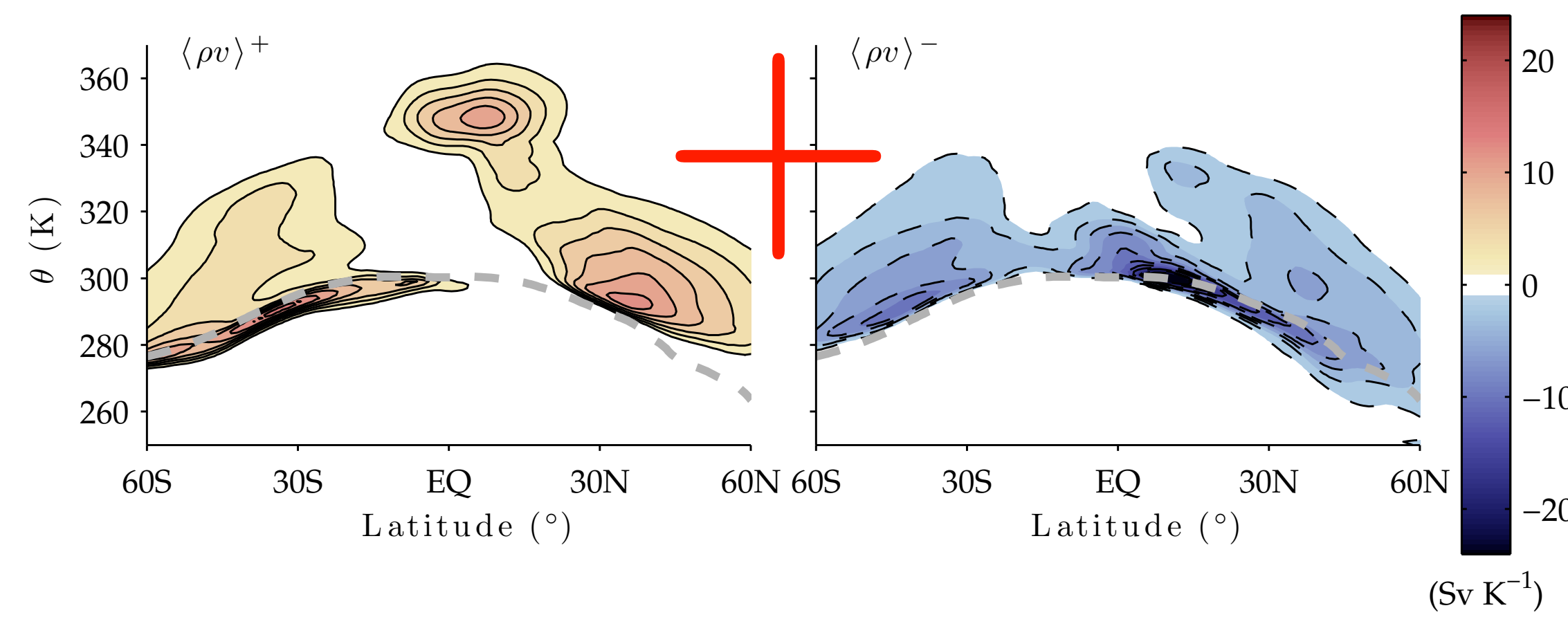**Mass Flux joint Distribution for DJF of 1981-2000 ERA 40 at 35N**

### 1) Directional Fluxes

Because poleward fluxes occupy a very different phase-space region than equatorward fluxes, it is convenient to consider each one in isolation. We project the positive and negative parts onto $\theta$:

$$\langle\rho v\rangle^+ = \frac{1}{2}\int_0^\infty (M+|M|)d\theta_e, \quad \langle\rho v\rangle^- = \frac{1}{2}\int_0^\infty (M-|M|)d\theta_e.$$

The sum of the two yields the meridional mass flux on $\theta$ surfaces, $\langle\rho v\rangle$:

$$\langle\rho v\rangle = \langle\rho v\rangle^+ + \langle\rho v\rangle^-.$$
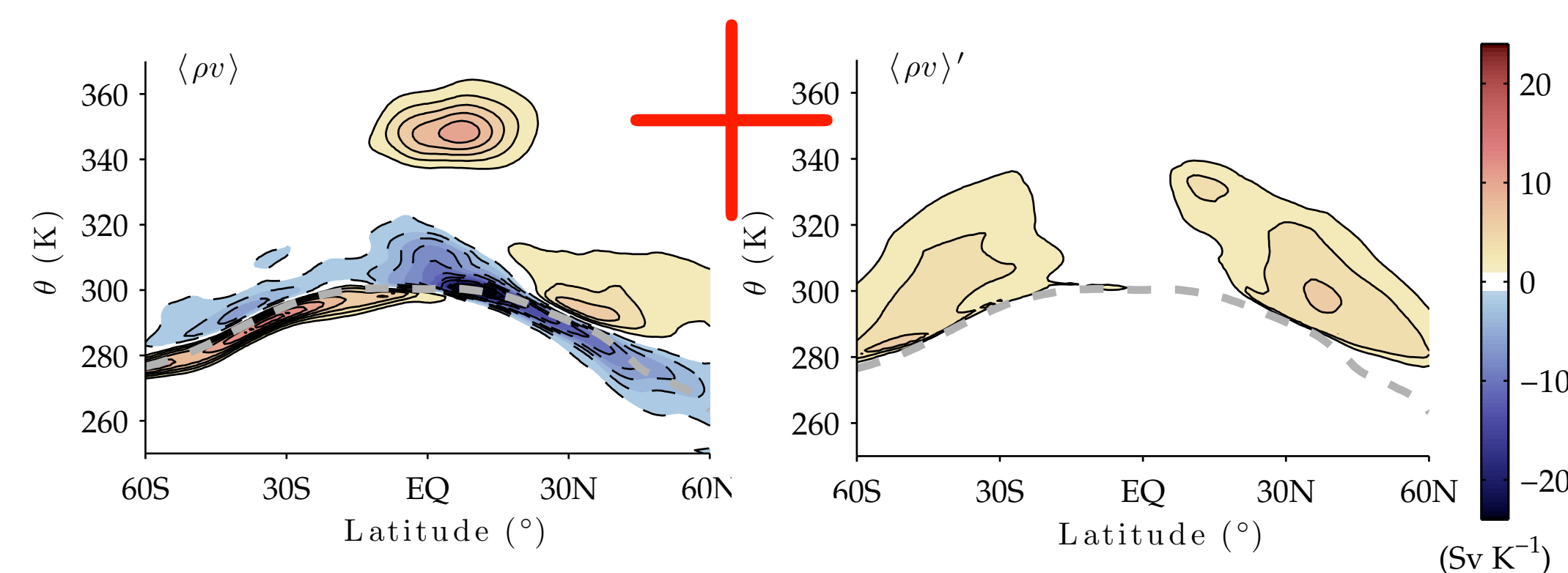


**Directional Fluxes for DJF of 1981-2000 ERA 40**

### 2) Moist Recirculation

This operation cancels out many mass fluxes. We call these cancellations the Moist Recirculation:

$$\langle\rho v\rangle' = \int_0^\infty |M|d\theta_e - \left|\int_0^\infty M d\theta_e\right|$$



**Net Mass Fluxes (left) and Moist Recirculation (right) for DJF of 1981-2000 ERA 40**
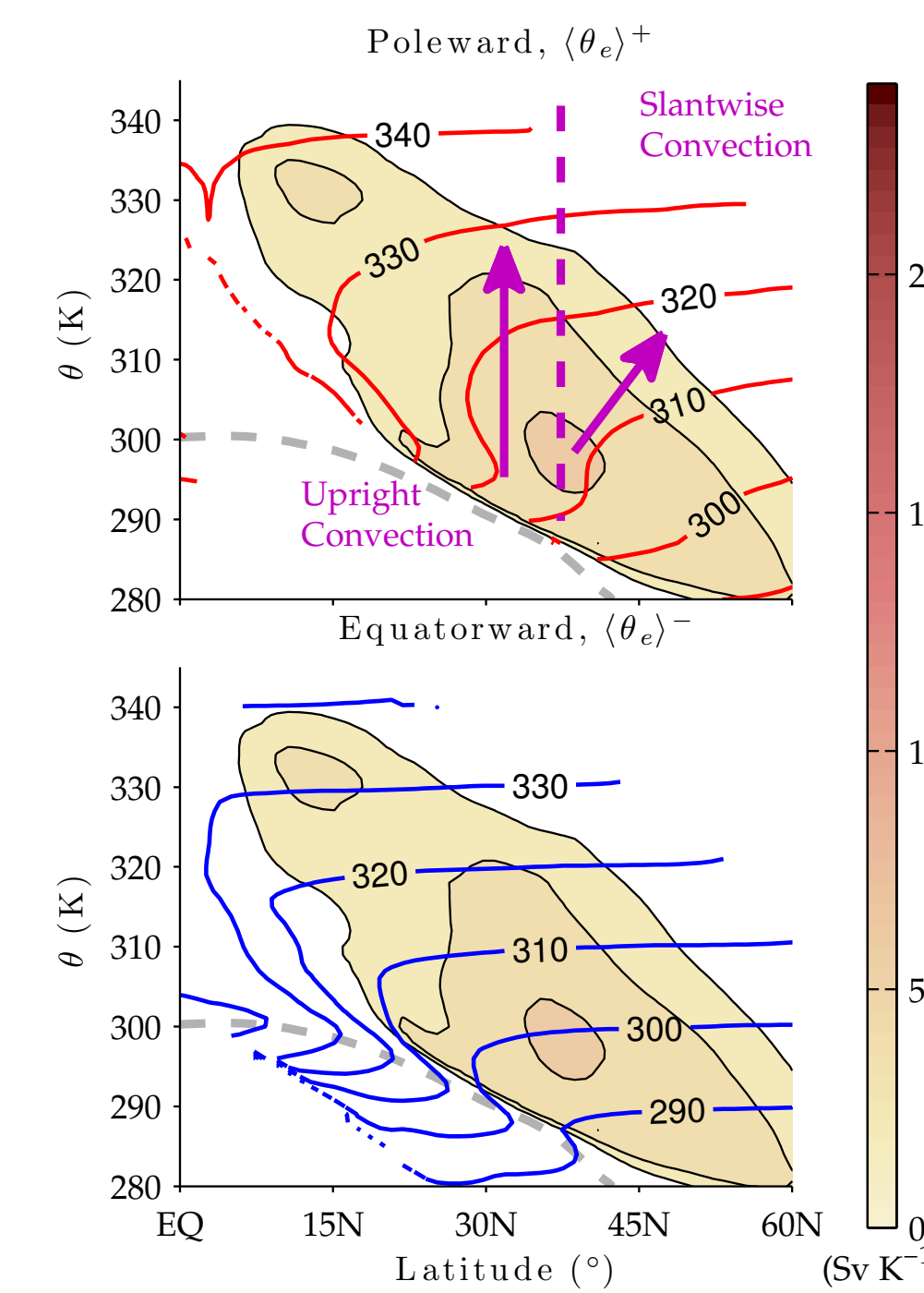
### 3) Interpretation[3]

Each of the directional components is associated with a specific $\theta_e$, their direction profile, that describe their moist isentropes:

$$\langle\theta_e\rangle^+ = \int_0^\infty \theta_e(M+|M|)d\theta_e \Big/ \int_0^\infty (M+|M|)d\theta_e,$$

$$\langle\theta_e\rangle^- = \int_0^\infty \theta_e(M-|M|)d\theta_e \Big/ \int_0^\infty (M-|M|)d\theta_e.$$

Poleward-moving moist parcels in the upper panel have a $\theta_e$ given by the red curves. Parcels on these curves have $\theta_e$ values that are much higher than their $\theta$ values. Equatorward-moving parcels in the lower panel have a $\theta_e$ given by the blue curves. They are dry in the mid troposphere and moist in the lower troposphere. The poleward $\langle\theta_e\rangle^+$ exhibits an ``S'' shape in the the subtropics but not in the midlatitudes. If a poleward-moving parcel in the subtropics was to follow its moist isentrope it would have to do so discontinuously: it would start along the lower part of the ``S'' and transition abruptly to the upper part. This is attributed to a predominance of fast upright convection. In the midlatitudes, no such ``S'' shape exists. Poleward-moving parcels can therefore follow their $\langle\theta_e\rangle^+$ moist isentrope all the way to the poleward edge of the recirculation. Such parcels ascend through the troposphere but do so slowly without the need for a fast transition to higher $\theta$. This is taken as an indication that slow slantwise convection is the dominant mode of convection in midlatitudes eddies.



**Directional $\theta_e$ and Moist Recirculation**

### 4) Computation

For the Mass Flux Joint distribution, the mathematical operations require meridional velocity, temperature and specific humidity at high-frequency (4x daily) computed on model levels for better accuracy. While its mathematical formulation is relatively simple, its computation from high-resolution data is expensive.

The output is a monthly and zonal mean quantity with a typical 10 fold reduction in size. Several other diagnostics, like EP fluxes, require a similar amount of input but have outputs several order of magnitudes smaller.

Computing $M(\phi, \theta', \theta'_e)$ is an $O(N)$ process: its processing time thus grows with the size of the dataset. In 2011, we processed the ERA-interim using a 2x Xeon E5603 quad core (48 Gb RAM) at U of T and in 2008 we processed the CMIP3 SRESA1B at NYU on the USQ cluster:



**Mass Flux Joint Distribution for DJF of 1991-2001 ERA Interim at 35N**

Processing time for some datasets

| Group | model | nlat, nlon, nlev | size (Gb) | Proc. Rate (Mb/s) | Transfer Rate (Mb/s) | year |
|---|---|---|---|---|---|---|
| CMIP5/RCP45 | MIROC4h | 320,640,56 | 7500 | ~7 | ? | 2011 |
| CMIP5/RCP45 | Can-ESM2 | 64,128,35 | 615 | ~7 | ? | 2011 |
| CMIP5/RCP45 | HadGEM2-E | 96,192,38 | 2010 | ~7 | ? | 2011 |
| ERA | Interim | 128,256,60 | 4500 | ~7 | ~6 | 2011 |
| CMIP3/SRESA1B | | ~T63,9 | ~1000 | ~0.4 | ~1.5 | 2008 |

This table illustrates how Moore's Law (processing time doubles every 18 months) overcomes Nielsen's Law (transfer time doubles every 24 months). While actual transfer rates make it possible to transfer the CMIP5 data at same rate as it is processed, in 10 years it will be ~3 times faster to process than to transfer the data.

## B. ExArch: Climate analytics on distributed exascale data archives

ExArch is a project funded by a G8 Research Councils Initiative on Multilateral Research. Its goal is to provide informatics solutions to the data storage and management problems facing the CMIP5/CORDEX archive and its potential successor. The project has three components: data management on the exascale, web service applications for server-side processing and quality assurance/advanced climate diagnostics.

The data management component will be done in collaboration with the Earth System Grid data servers. It will seek ways to better incorporate the METAFOR questionnaire and its expanded version with the data. These steps should allow fully automatized data retrieval for the efficient intercomparison of future model data.

UNIVERSITY OF TORONTO    British Atmospheric Data Centre    DKRZ    Centro Euro-Mediterraneo per i Cambiamenti Climatici    UCLA    Princeton University    Institut Pierre Simon Laplace
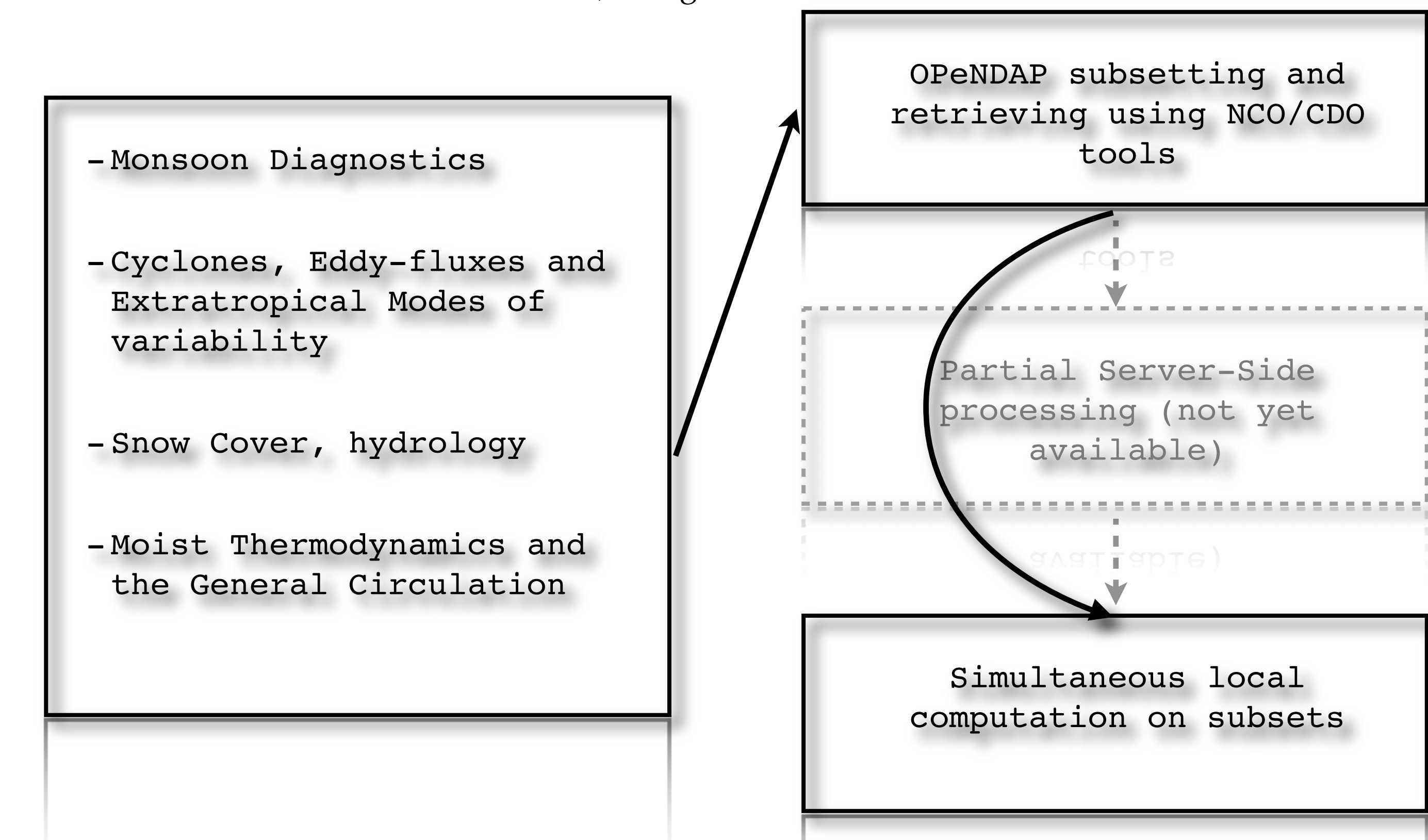
## C. Web and Server-side Processing

One of the main outcomes of the ExArch project will be the development of a web processing service that will enable several operations to be performed by the data server, before the data transfer. Performing data-reducing operations first will reduce the transfer and the storage requirement by the climate scientist and thus mitigate issues arising from the handling of large datasets. The implementation will be based on Climate Data Operators (CDO) with available operations being added one by one over the next 3 years.

## D. Climate Diagnostic Benchmark (CDB)

The University of Toronto team is responsible for the development of a suite of climate diagnostics to take full advantage of the web processing framework. The CDB will evolve as more features will become available in server-side computations.

The first versions of the CDB will implement the Advanced Climate Diagnostics (e.g. Mass Flux Joint Distribution) using an OPeNDAP framework:

- Monsoon Diagnostics
- Cyclones, Eddy-fluxes and Extratropical Modes of variability
- Snow Cover, hydrology
- Moist Thermodynamics and the General Circulation

OPeNDAP subsetting and retrieving using NCO/CDO tools

Partial Server-Side processing (not yet available)

Simultaneous local computation on subsets

## References

1. ExArch Delivery Plan (WP1), 2011: http://proj.badc.rl.ac.uk/exarch/wiki/DeliveryPlan
2. Pauluis, O., A. Czaja and R. Korty (2008), The global circulation on moist isentropes. *Science*, 321, 1075-1078.
3. Laliberté, F., T. Shaw and O. Pauluis (2011), Moist recirculation and water vapor transport on dry isentropes. To appear in *JAS*.