

## Modélisation numérique du climat sur architectures haute performance : quelques succès et défis à l'IPSL

Thomas Dubos LMD/IPSL, École Polytechnique

Marie-Alice Foujols

IPSL, CNRS

Yann Meurdesoif

LSCE/IPSL, CEA

## Modélisation numérique du climat sur architectures haute performance : quelques succès et défis à l'IPSL

- Contexte et quelques chiffres
- Du vectoriel au massivement parallèle
- Vers l'exascale ? des défis techniques, mais pas seulement

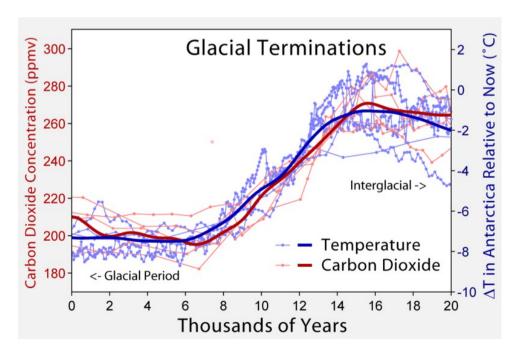
Le temps : une réalisation

- Simuler
- Prévoir

Le climat : une statistique et son évolution

- Modéliser
- Comprendre/expliquer
- Anticiper





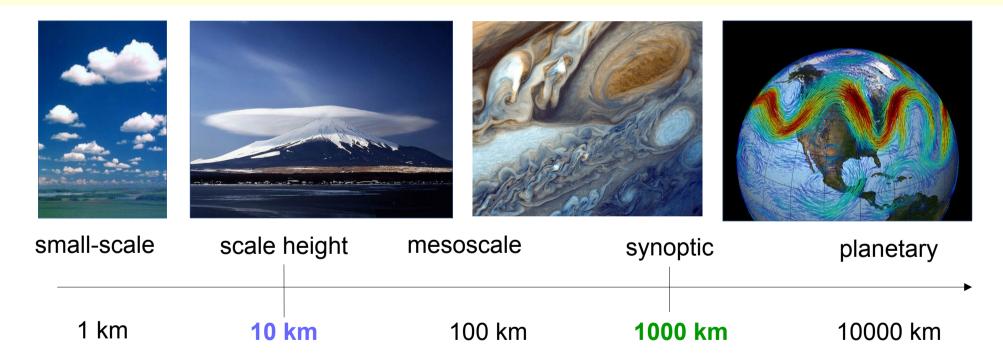
## The atmosphere : a gravity-dominated, compressible flow

#### Characteristic scales

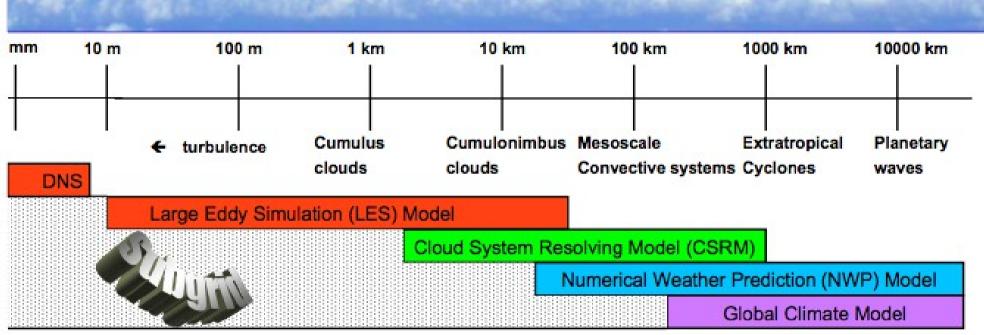
Velocity: Sound c ~ 340m/s Wind U ~ 30m/s

• Time: Buoyancy oscillations  $N \sim g/c \sim 10^{-2} s^{-1}$  Coriolis  $f \sim 10^{-4} s^{-1}$ 

• Length: Scale height H=c²/g=10km Rossby radius: R=c/f ~ 1000 km

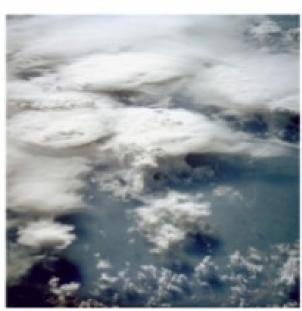


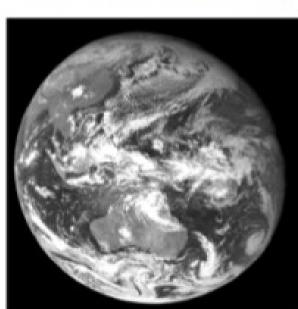
## Les nuages : différentes échelles, différents processus









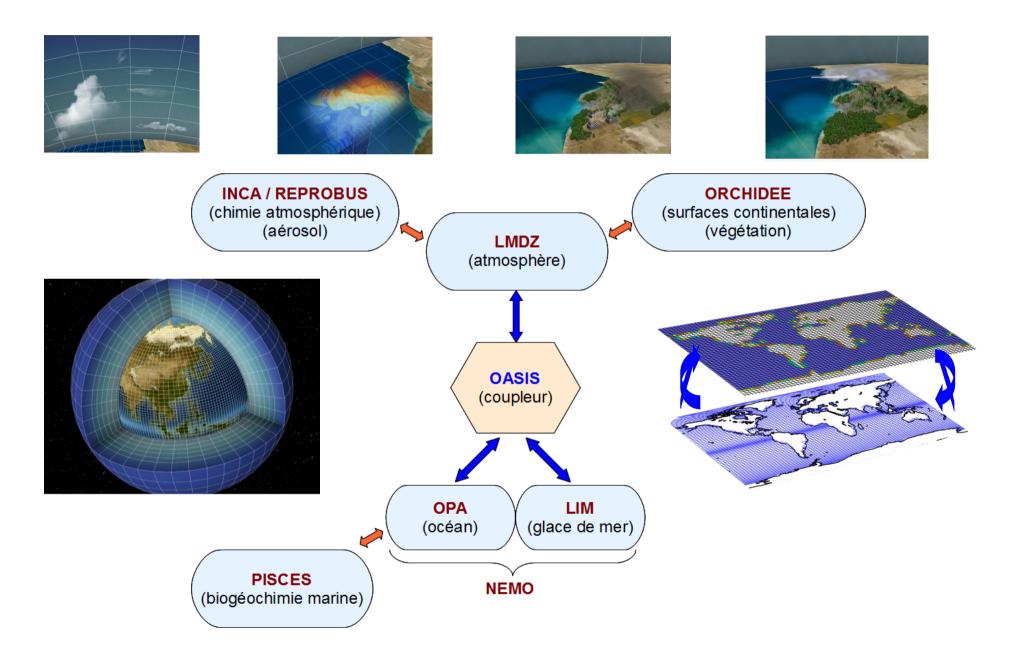


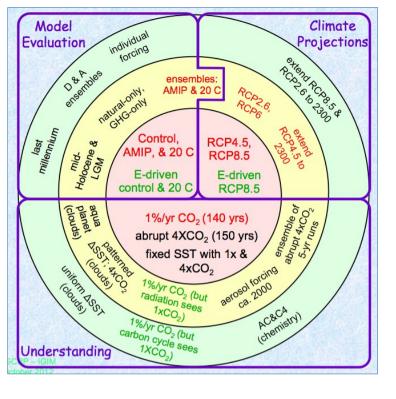
P. Siebesma

## La modélisation des nuages convectifs

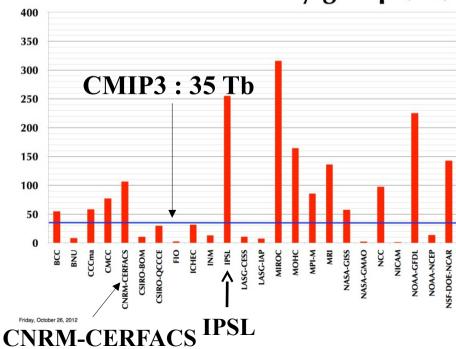
Modèle de circulation générale (GCM): dx=dy=de 50 à 400km Ζ dz= de 100m à 2km 10km 3km 2km 1km 100m MO MO Q X C. Rio

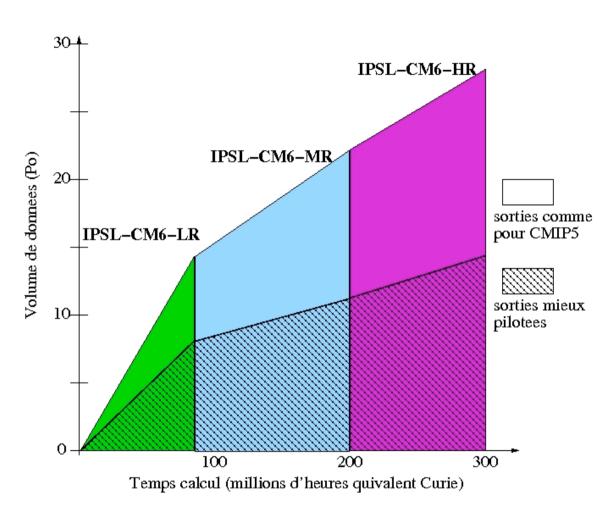
## Earth System Modelling at IPSL





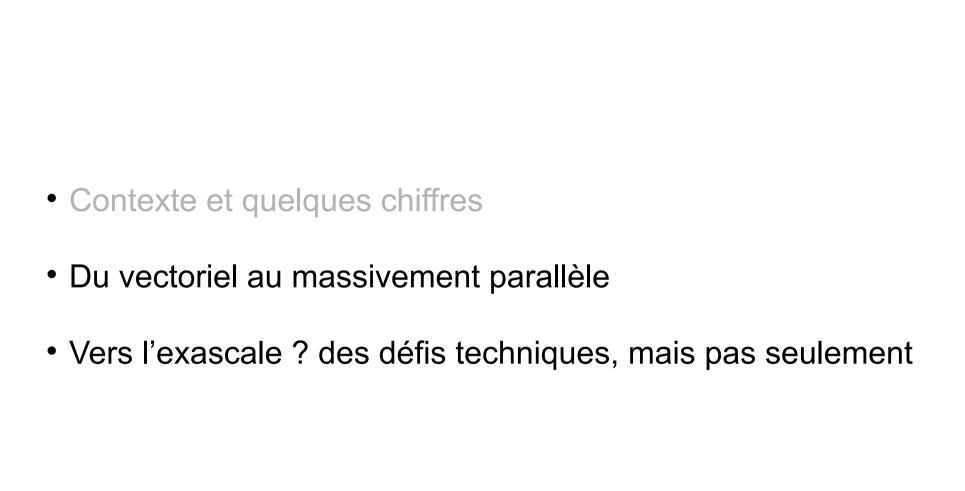
## **CMIP5** data volumes by group (TB)



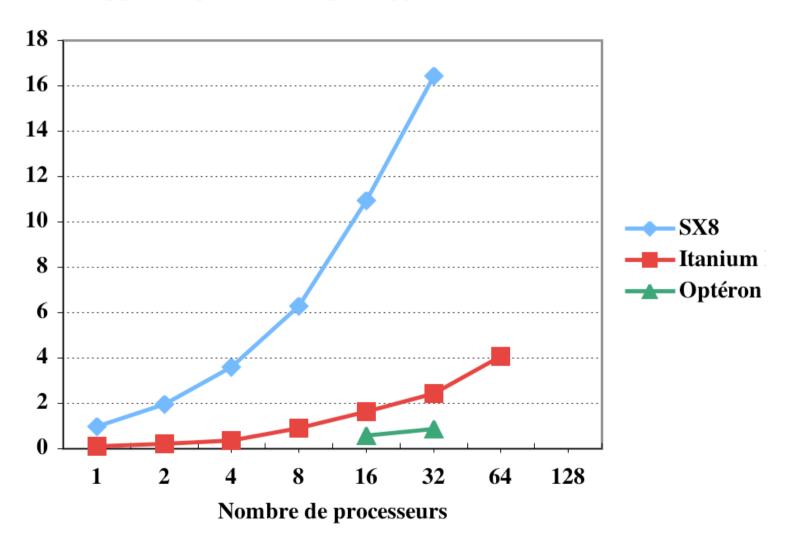


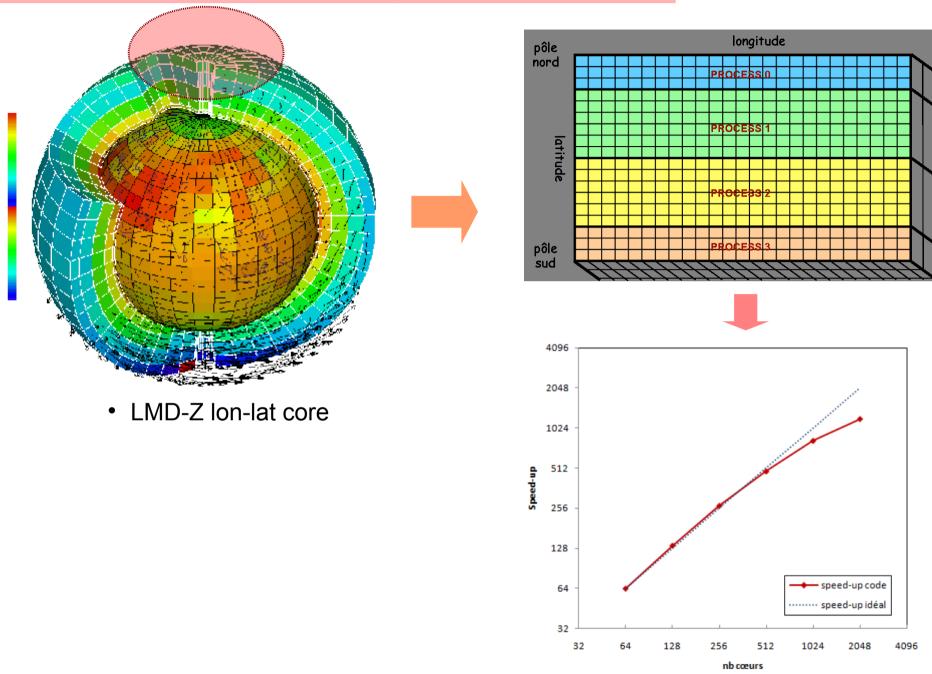
## 10-year perspective for CMIP

	CMIP5	CMIP6	CMIP7
Year	2012	2017	2022
Power factor	1	30	1000
Npp	200	357	647
Resolution [km]	100	56	31
Number of mesh points [millions]	3,2	18,1	108,4
Ensemble size	120	214	388
Number of variables	800	1068	1439
Interval of 3-dimensional output (hours)	6	4	3
Years simulated	90000	120170	161898
Storage density	0,00002	0,00002	0,00002
Distributed Archive Size (Pb)	3,19	86,05	2260,20



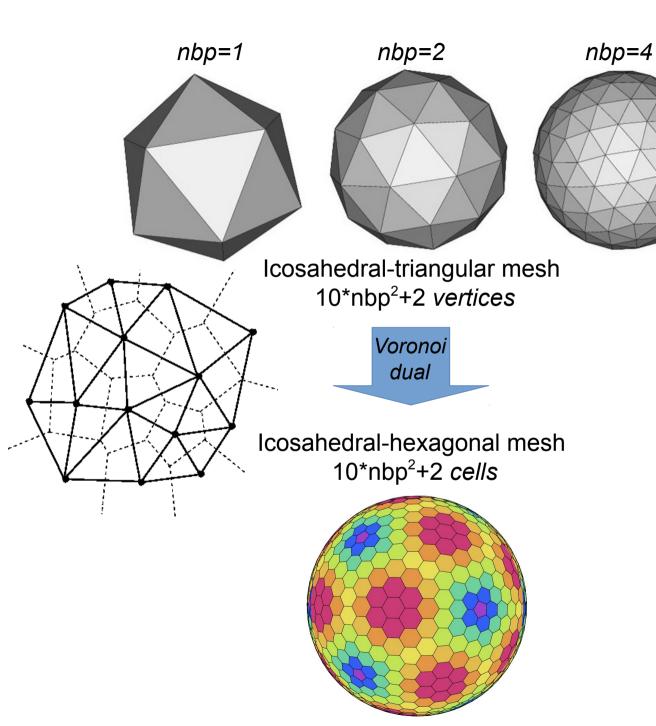
#### Rapport de performance par rapport à 1 CPU SX8



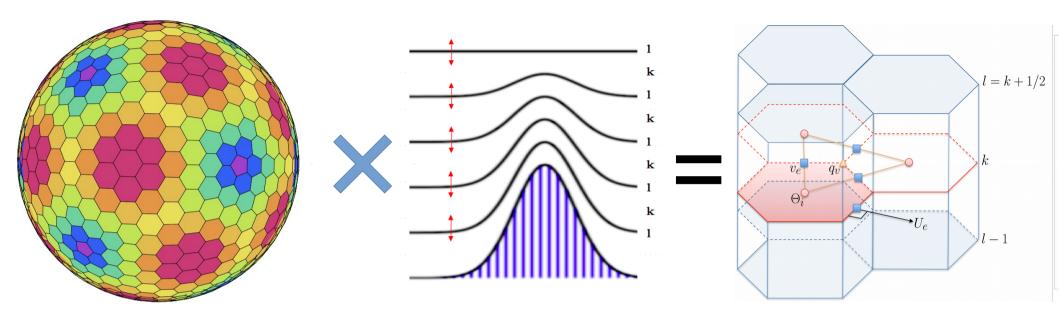


Y. Meurdesoif, 2010 1/4 degré (25km) => max. 240 processus MPI

## Mesh partitioning for parallel computing



- Easy to partition into 10 x nsplit<sup>2</sup> domains
- About (nbp/nsplit)<sup>2</sup> cells per domain = MPI process
- Nbp/nsplit>10 for performance
- Ex: 25km ~10<sup>6</sup> cells ~10 000
   MPI processes

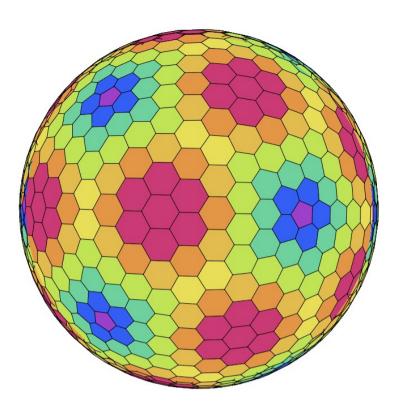


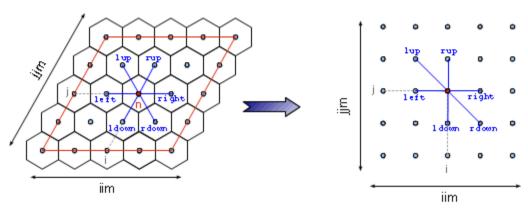
- Discrete integration by parts (Bonaventura & Ringler, 2005; Taylor, 2010)
- Energy- and vorticity- conserving Coriolis discretization
- (TRiSK: Thuburn et al., 2009; Ringler et al., 2010)



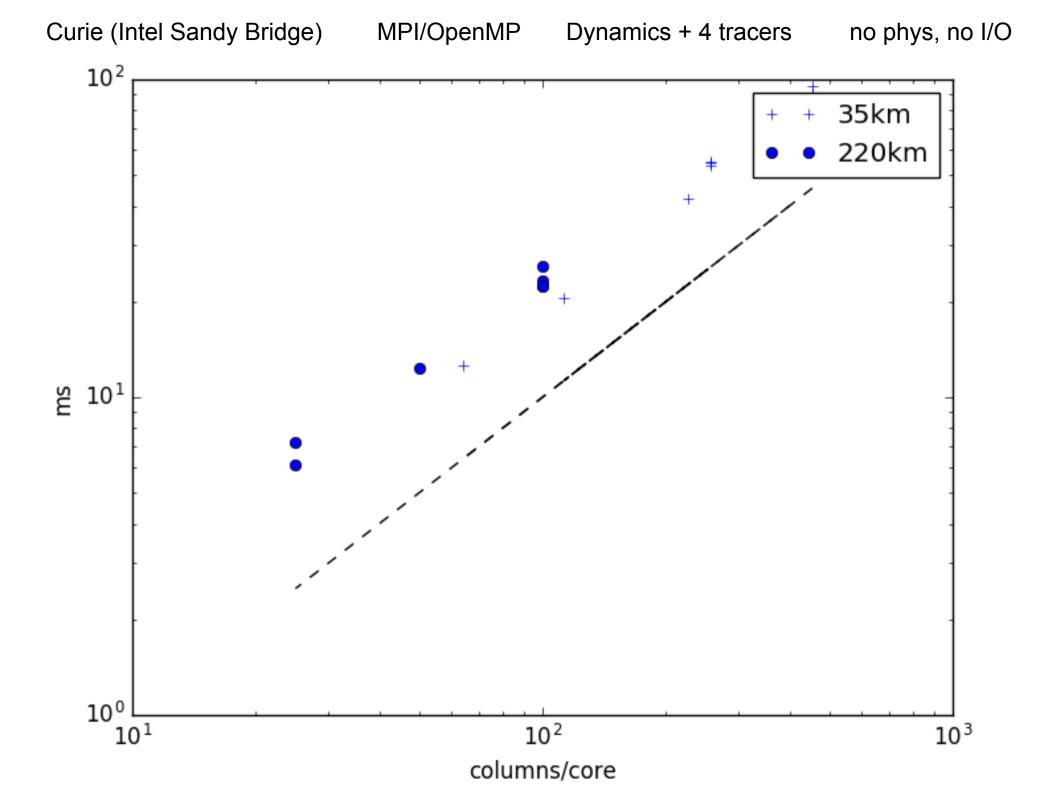
# Energy-conserving 3D core

(Tort & Dubos, 2015;
 Dubos et al., 2015)





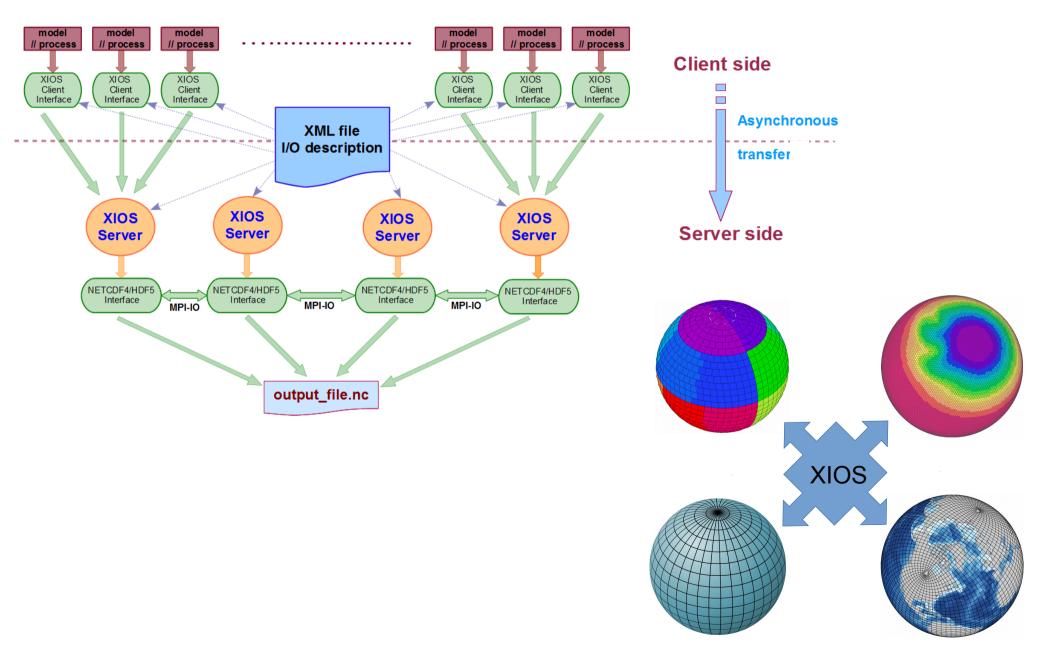
- Vertical direction in outer loops
- Direct access to neighbours via constant offsets
- No special case for pentagons (handled by metrics)



- CMIP requires a throughput of x10000 (30SYPD)
- Some climate modelling still doable with x1000 (3SYPD)
- Ability to attain x1000 depends on maximum stable time step (numerics) and walltime needed to perform one time step (implementation)
- Assuming a large enough machine, reducing walltime is a strong scaling problem
- For DYNAMICO, dt (in sec) is about 2.5\*dx (in km)
- => 3SYPD
  - At 25km resolution requires about 60 ms per full time step
  - At 8km resolution requires about 20 ms per full time step
  - At 1km resolution requires about 2.5 ms per full time step

## XIOS (Y. MEURDESOIF): XML I/O SERVER

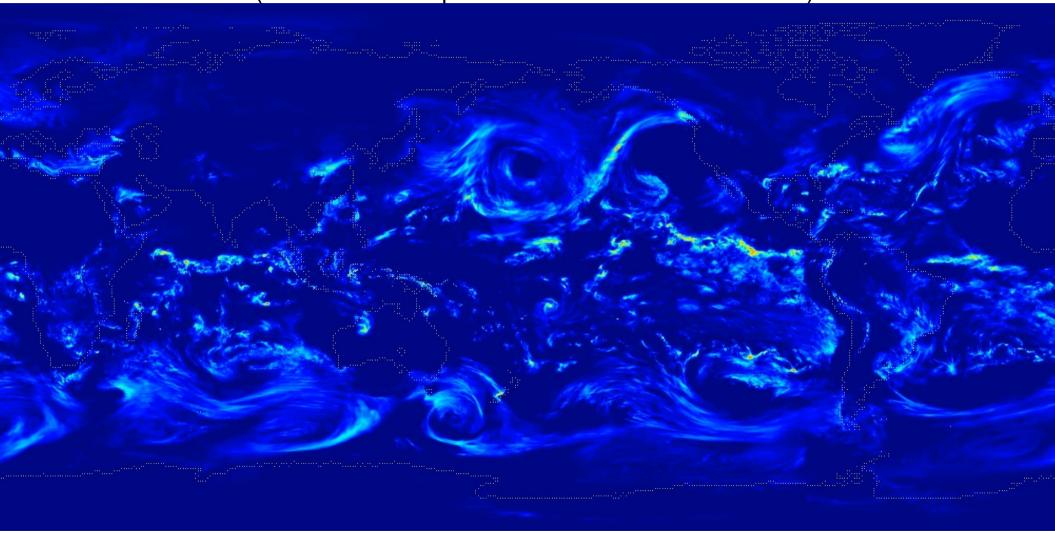
# PARALLEL ASYNCHRONOUS I/O - ONLINE POST-PROCESSING LIBRARY AND SERVER



## IPSL-CM7A-HR:

DYNAMICO-LMDZ, 25km

(1 024 000 atmospheric columns x 80 vertical levels)

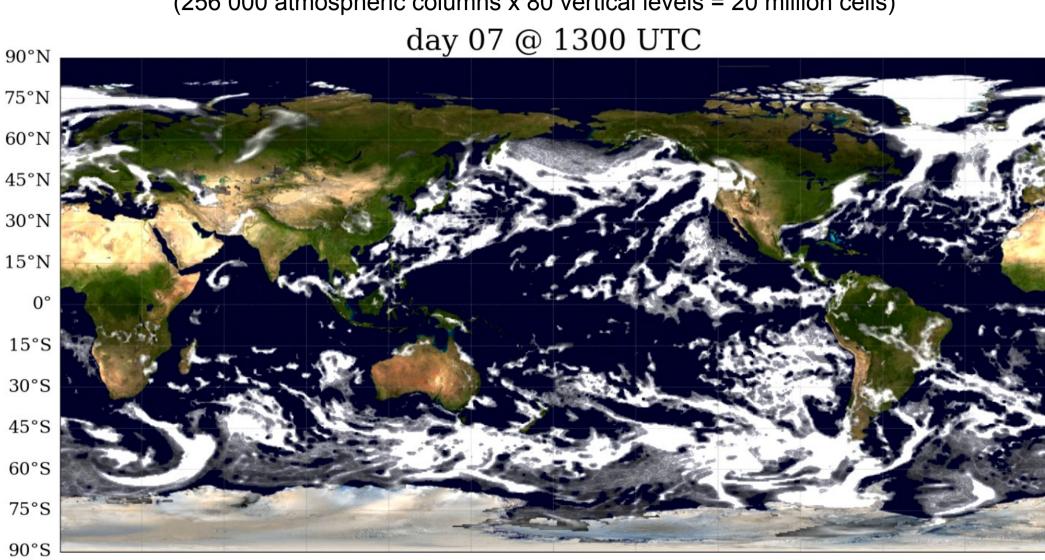


Liquid water

## IPSL-CM7A-HR:

DYNAMICO-LMDZ, 50km

(256 000 atmospheric columns x 80 vertical levels = 20 million cells)



Low-level cloudiness

180°

150°W

120°W

90°W

60°W

30°W

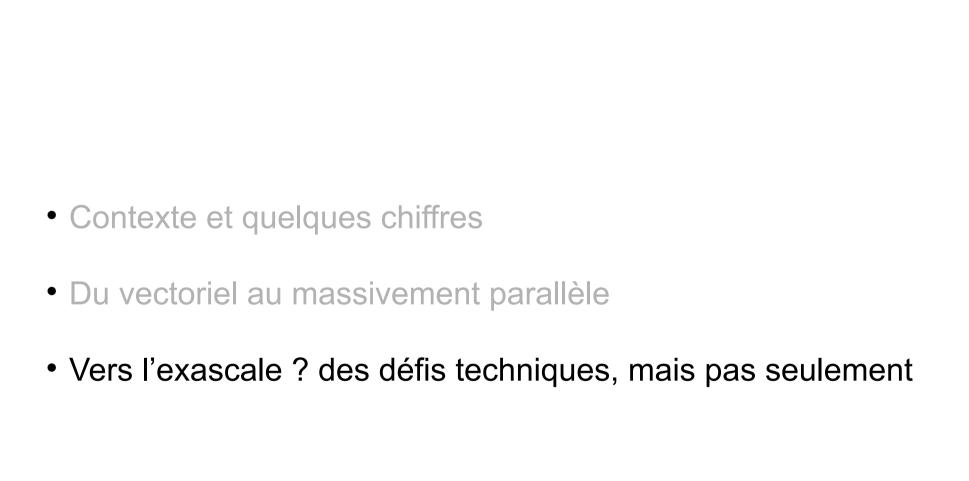
150°E

30°E

60°E

90°E

120°E



#### Défi technique : adaptation des codes à de nouvelles architectures

Ex : portage de DYNAMICO sur GPU Nvidia V100 (Y. Meurdesoif, A. Durocher (LSCE/IPSL) / HPE / contrat de progrès IDRIS)

« Le contrat de progrès a pour objectif d'accompagner les communautés utilisatrices de l'IDRIS ... le but final étant, en cas de succès, de pouvoir migrer ces communautés d'utilisateurs identifiées sur la partition convergée... »

« Pour être validée, la version GPU d'un code devra ... être en moyenne sur l'ensemble des cas tests au moins 4 fois plus performante que les versions non-accélérées en comparaison nœud à nœud. »

DYNAMICO: ~10 000 LOC dont ~2000 LOC critiques

=> portage manuel

- par insertion de directives OpenACC
- adaptation de la gestion mémoire, halos, ...

Nb colonnes atmosphère	16000	64000	256000
40 procs MPI (CPU)	68 ms	282 ms	1193 ms
1 GPU	26,6 ms	77 ms	299 ms
Accélération	2,5	3,6	4,0

```
!$acc parallel loop collapse(2)
DO 1 = ll_begin, ll_end
DO ij = ij_begin_ext,ij_end_ext
    uu_right = 0.5*(rhodz(ij,1)+rhodz(ij+t_right,1))*u(ij+u_right,1)
    uu_right = uu_right*le_de(ij+u_right)
    hflux(ij+u_right,1) = uu_right
    uu_lup = 0.5*(rhodz(ij,1)+rhodz(ij+t_lup,1))*u(ij+u_lup,1)
    uu_lup = uu_lup *le_de(ij+u_lup)
    hflux(ij+u_lup,1) = uu_lup
    uu_ldown = 0.5*(rhodz(ij,1)+rhodz(ij+t_ldown,1))*u(ij+u_ldown,1)
    uu_ldown = uu_ldown*le_de(ij+u_ldown)
    hflux(ij+u_ldown,1) = uu_ldown
END DO
```

## Défi technique : adaptation des codes à de nouvelles architectures

LMDZ (paramétrisations sous-maille): ~100 000 LOC

- Profil très plat, pas de 'hot spot'
- Rétrocompatibilité : on supprime rarement du code existant => parties très anciennes (>30 ans)
- Code très touffu
- Parties écrites « maison » : évoluent en permanence
- Parties importées (code radiatif RRTM) : expertise manquante

Si des moyens étaient obtenus pour le faire :

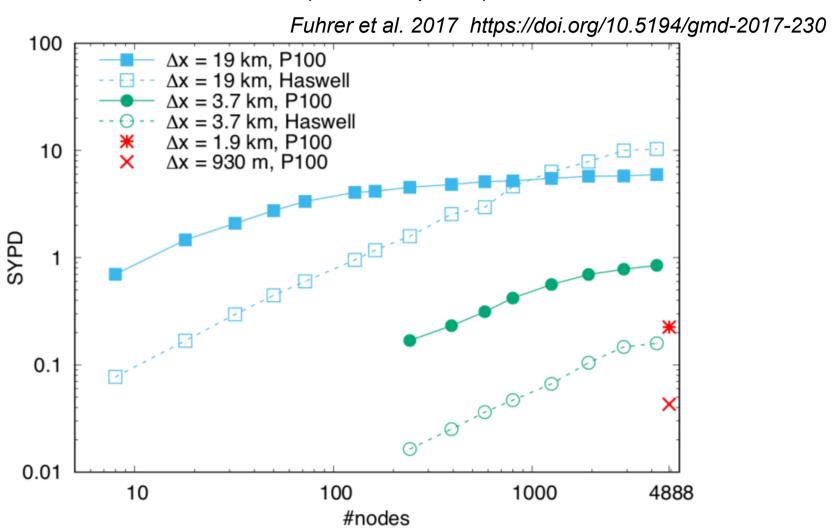
- Est-ce un investissement (architecture pérenne) ou une dépense (architecture sans lendemain) ?
- Comment être performant sur plusieurs architectures avec un code unique ?
- Comment faire évoluer après portage un code communautaire écrit par des physiciens ?

=> questions liées au développement communautaire des codes Dans l'air du temps :

- 'separation of concerns', 'domain-specific languages'
- machine learning

## Dans l'air du temps : 'separation of concerns' et domain-specific languages

Modèle COSMO + DSL Stella/GridTools (C++ / templates)

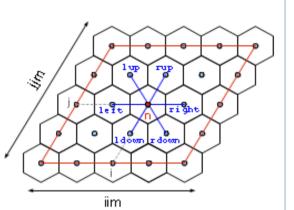


Écueils potentiels : pérennité du DSL, « learning curve », usine à gaz ?

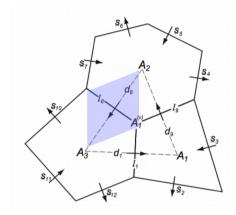
Dans l'air du temps : 'separation of concerns' et domain-specific languages

## DySL: DYNAMICO-specific « language »

An experiment in code maintainability and separation of concerns



```
KERNEL('div')
50
     FORALL CELLS EXT()
       ON PRIMAL
51
         div ij=0.
52
         FORALL EDGES
53
           div_ij = div_ij + SIGN*LE DE*u(EDGE)
54
55
         END BLOCK
         divu(CELL) = div ij / AI
56
       END BLOCK
57
58
     END BLOCK
59 END BLOCK
```

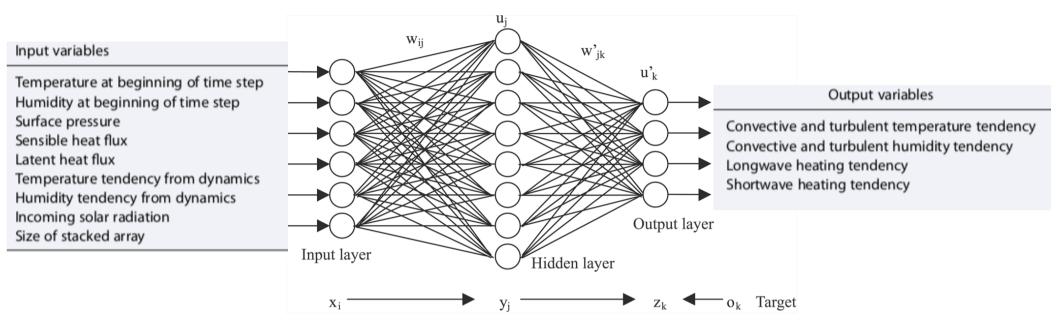


```
DO 1 = 11 begin, 11 end
         DO ij=ij begin ext, ij end ext
            div ij=0.
            div ij = div ij + ne rup*le de(ij+u rup)*u(ij+u rup,l)
            div ij = div ij + ne lup*le de(ij+u lup)*u(ij+u lup,1)
 7
            div_ij = div_ij + ne_left*le_de(ij+u_left)*u(ij+u_left,l)
 8
            div ij = div_ij + ne_ldown*le_de(ij+u_ldown)*u(ij+u_ldown,1)
            div_ij = div_ij + ne_rdown*le_de(ij+u_rdown)*u(ij+u_rdown,1)
10
            div ij = div ij + ne right*le de(ij+u right)*u(ij+u right,l)
11
            divu(ij,l) = div ij / Ai(ij)
12
         END DO
13
      END DO
14
```



- Simplicité plutôt que généralité => faible coût d'entrée
- Génération de blocs de code Fortran par manipulation du pseudo-code « DySL »
- => faible coût de sortie
- Approches plus sophistiquées (parser/backend) envisagées avec J. Bigot (MdLS)

#### Dans l'air du temps : machine learning



Stephan Rasp<sup>a,b,1</sup>, Michael S. Pritchard<sup>b</sup>, and Pierre Gentine<sup>c,d</sup> www.pnas.org/cgi/doi/10.1073/pnas.1810286115

- Appel à routine(s) complexes/obscures remplacé par réseau de neurone
- => générique, potentiellement efficace sur GPU/architecture convergée
- Encore beaucoup de questions : contraintes physiques, utilisation hors domaine d'apprentissage ...

- IPSL : une transition vers le parallélisme massif opérée et opérationnelle (CMIP6)
- Vers l'exascale ? des défis techniques, mais pas seulement
  - Exploiter le parallélisme massif peut nécessiter une reformulation profonde du contenu du modèle (ex : DYNAMICO)
  - Les architectures exascale auront-elle la scalabilité forte nécessaire pour la modélisation du climat ?
  - Le but n'est pas de porter nos codes, mais d'arriver à un modèle de développement qui assure une certaine portabilité de la performance
  - Les solutions pour parvenir à ce but doivent tenir compte du caractère communautaire du développement, et du rôle de laboratoire scientifique joué par le modèle
    - solutions « maison » : sur mesure mais ressources, expertise ?
    - solutions externes : besoins spécifiques, pérennité, coût de sortie, learning curve ?
    - contourner le problème : machine learning ??

#### Merci de votre attention

