

Retour d'expérience sur BlueGene : CERFACS, EDF, CEA



- Saturne (CFD RANS) --> 1024 proc
- Zephyr (Burgers 2D) --> 1024 proc
- Dymoka (Dynamique Moléculaire) --> 1024 proc
- VASP (ab initio DFT) --> 32 proc



- TRIO U (CFD) --> 2048 proc



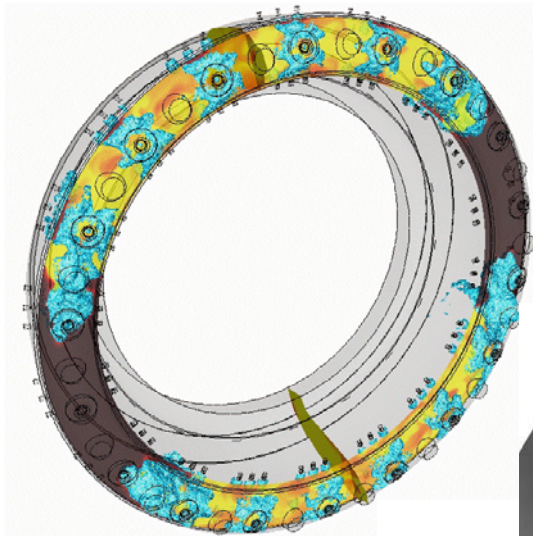
- AVBP (CFD LES combustion) --> 5120 proc
=> BlueGene + ter@10



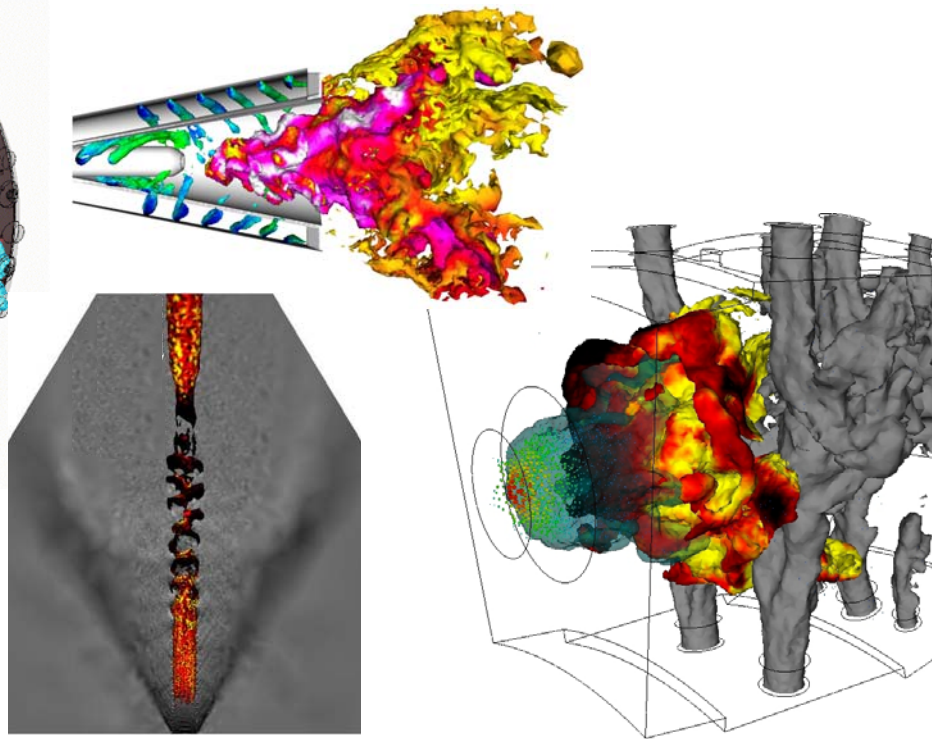
Simulations de combustion turbulente avec AVBP sur calculateurs massivement parallèles

Y. Sommerer, M. Boileau, G. Staffelbach, T. Schonfeld, T. Poinsot, L. Giraud

CFD Team - CERFACS



www.cerfacs.fr/cfd



AVBP est développé au CERFACS depuis 10 ans pour simuler les instabilités de combustion qui ont pour sources des phénomènes hydrodynamiques et/ou acoustiques.

LES réactif compressible (explicite) et diphasique sur maillages hybrides construit dès son origine pour les machines massivement // (mpi)

C'est un code développé par une grosse équipe de chercheurs multi-sites:

CERFACS (30), IFP (10), EM2C (10), IMFT (3)

AVBP: un code de recherche géré comme un code industriel

De par son avancée scientifique, AVBP a peu de concurrents au niveau international dans le domaine de la combustion turbulente numérique.

AVBP est choisi par de nombreux laboratoires de recherche et d'industriels et est utilisé dans de nombreux projets Européens:

- EM2C, IMFT, IRPHE, Univ. Montpellier, Univ. Jussieu (lab. Jacques-Louis Lions), Univ. Belfast, Univ. Zaragossa, Univ. Twentee...
- Snecma, Turbomeca, Siemens, Alstom, EDF, GDF, Air Liquide, Peugeot, Ferrari, ONERA, IFP...
- PRECCINSTA, MOLECULES, INTELLECT, FUEL CHIEF, DESIRE, LESSCO2, FLUITSCOM, TLC...

On a besoin de gros moyens CPU car:

- temps physiques à simuler longs
- besoin de multiphysique (diphasique, couplage avec codes rayonnement et thermique etc)
- extension du domaine de calcul (direction axiale : x 1.2 à 1.5
direction azimutale : x 20)

Taille maillages max. en 2004 : 3 à 5 millions de cellules

=> on doit passer à 50 ou 100 millions!!!

- besoin de diminuer le temps de restitution

Besoin d'un saut technologique, c'est ce qu'apporte les nouvelles machines massivement parallèles.

Tests de scalabilité: le but est de montrer la faisabilité d'une LES sur chambre complète avec des temps de restitution acceptables.

3 configurations:

Chambre Siemens: 24 bruleurs ($40 \cdot 10^6$ cells)

Chambre Turboméca: 18 bruleurs ($18 \cdot 10^6$ cells)

- injection carburant gazeux

- injection carburant liquide

2 machines:

IBM BlueGene/L (jusqu'à 5120 processeurs)

Ter@10 (jusqu'à 1900 processeurs)

Configuration 1: chambre Siemens 24 bruleurs 40 millions de cellules

=> 1^{ère} mondiale: c'est encore à l'heure actuelle la plus grosse simulation jamais calculée en combustion turbulente LES !

Collaboration IBM - CERFACS:

Accès aux BlueGene/L de Rochester et Thomas Watson => fenêtre de run 1 mois

Aucun pb de portage (sauf alloc. mémoire excessive de MPI_PROBE)

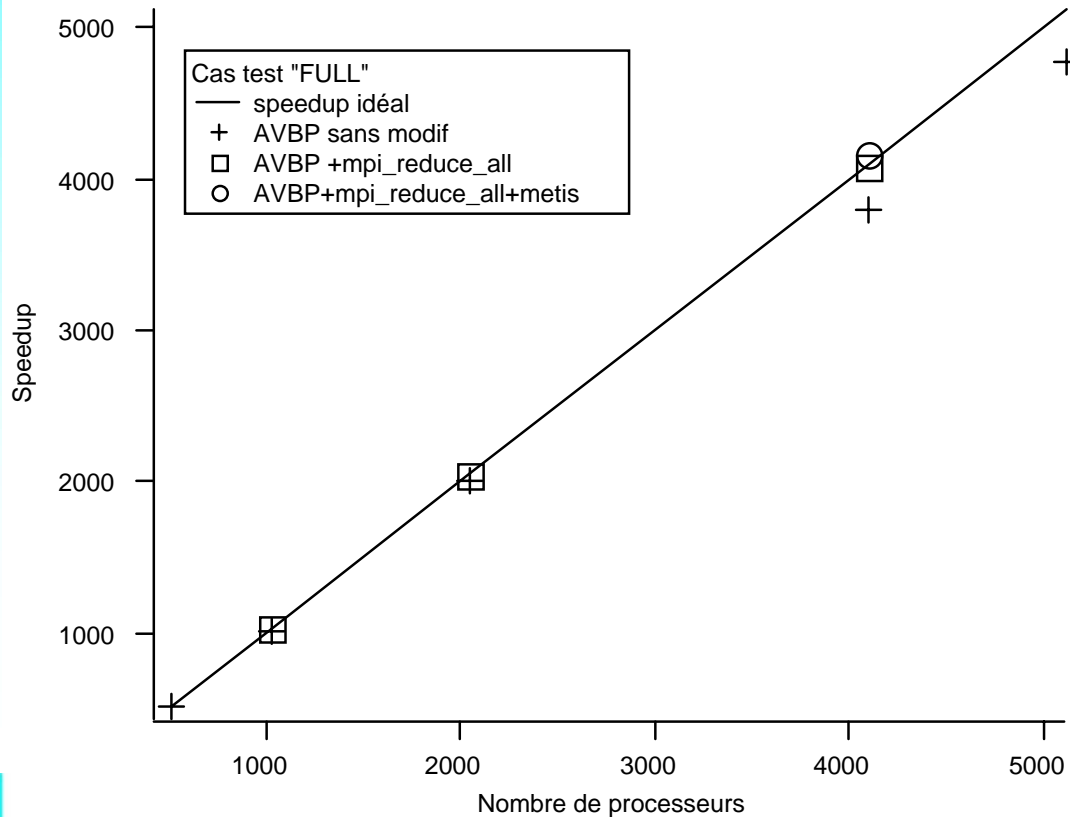
QuickTime™ et un
décompresseur TIFF (LZW)
sont requis pour visionner cette image.

QuickTime™ et un
décompresseur TIFF (LZW)
sont requis pour visionner cette image.

QuickTime™ et un
décompresseur TIFF (LZW)
sont requis pour visionner cette image.

QuickTime™ et un
décompresseur TIFF (LZW)
sont requis pour visionner cette image.

Les performances: "FULL" 40M cells

QualiTech™ et les
autres marques "IP" (TM)
sont registrés pour identifier cette page.

On compare 3 séries :

- 1/ sans modif
- 2/ optim réductions collectives
- 3/ optim découpage

Plusieurs remarques:

1/ Perfs excellentes!!!!

Speedup de **4078 sur 4096** procs (coût // 0.4%) et 4770 sur **5120**...
beaucoup de codes CFD aimeraient atteindre des speedup de 10!!!

2/ Avec réductions collectives on exploite mieux les perfs des réseaux d'interconnection de BlueGene => gain 7,5%

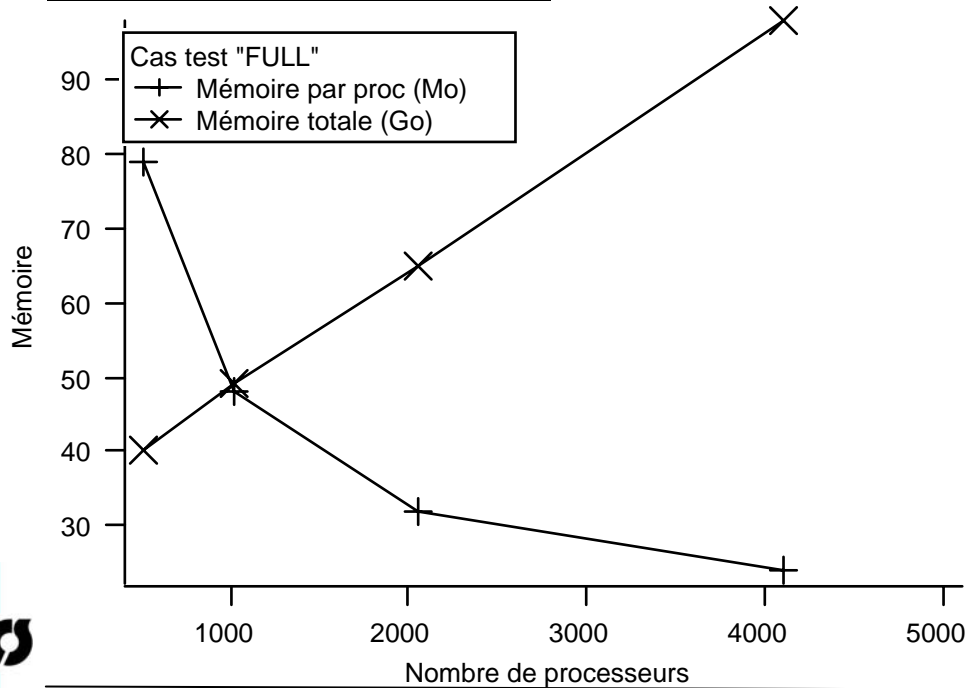
3/ Gain CPU dû à l'optimisation du découpage de maillage de 2%

Quintessence™ est un
enregistrement TM de
son propriétaire. TM © 2001
tous droits réservés pour ce site.

Traces MPI

QuickTime™ et un
décompresseur TIFF (LZW)
sont requis pour visionner cette image.

Ressources mémoire



Nombre de processeurs	Mémoire par processeur	Mémoire totale
512	79 Mo	40 Go
1024	48 Mo	49 Go
2048	32 Mo	65 Go
4096	24 Mo	98 Go

Les performances: "FULL" 40M cells

Qualitas™ et
distributeur TFF (2001)
ont reçu pour obtenir cette page.

Virtual node mode vs Coprocessor mode

En Coprocessor mode (Co), 1 seul des 2 procs du nœud travaille mais il a accès à la totalité des 512 Mo de mémoire (l'autre proc gère les accès mémoire et messages MPI).

En Virtual node mode (Vn) mode, les 2 processeurs du nœuds travaillent mais il ont chacun accès à 256 Mo de mémoire.

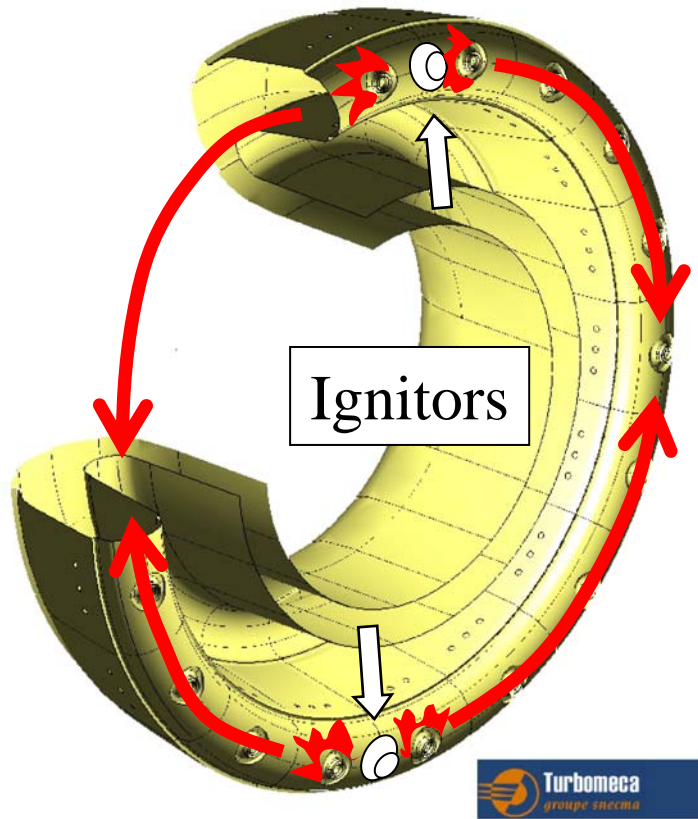
Pour un nombre de nœuds identique, Vn devrait être 2 fois plus rapide que Co.

Pour AVBP:

$$\text{CPU Co} = 1.8 \text{ CPU Vn}$$

Donc en Vn on a des conflits mémoire (10% CPU)... mais qui restent faibles (<50%) => on gagne en temps de restitution en Vn à nombre de nœud identique.

Configuration 2: chambre Turbomeca 18 bruleurs, 18 millions cellules
=> 1^{ère} mondiale: c'est le seul calcul d'allumage de chambre annulaire
complète !



Collaboration CERFACS-TURBOMECA-
IBM et CERFACS-CINES:

- * Accès aux BlueGene/L de Rochester et Thomas Watson (1024 a 2048 processeurs)
- * Accès aux machines du CINES (SGI 03800 128 processeurs)

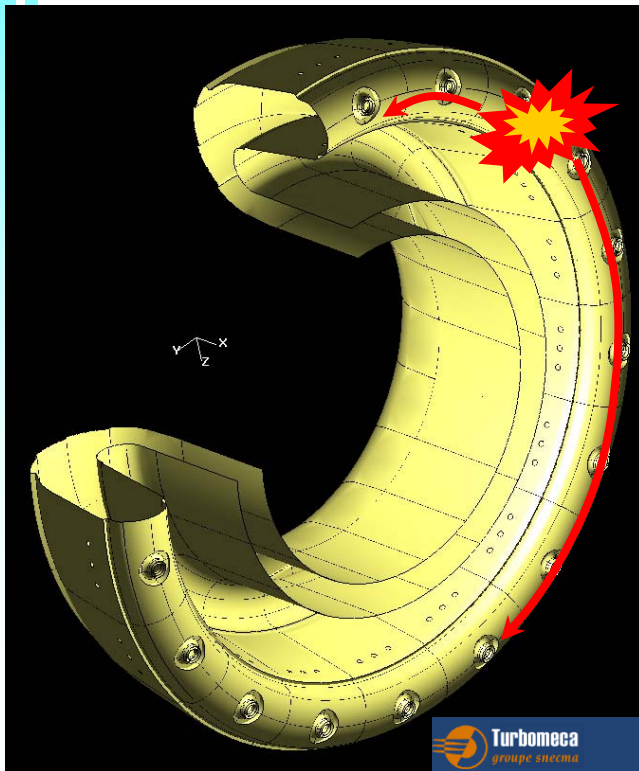
Comment se propage l'allumage d'un
brûleur à l'autre?
Besoin de calcul le foyer annulaire
complet (au moins plusieurs brûleurs)

Résultat:

QuickTime™ et un
décompresseur Photo - JPEG
sont requis pour visionner cette image.

Configuration 3: chambre Turbomeca 18 bruleurs, 18 millions cellules INJECTION CARBURANT LIQUIDE

=> 1^{ère} mondiale: c'est le seul calcul d'allumage de chambre annulaire complète avec injection de carburant liquide !



Collaboration CERFACS-CEA-TURBOMECA:

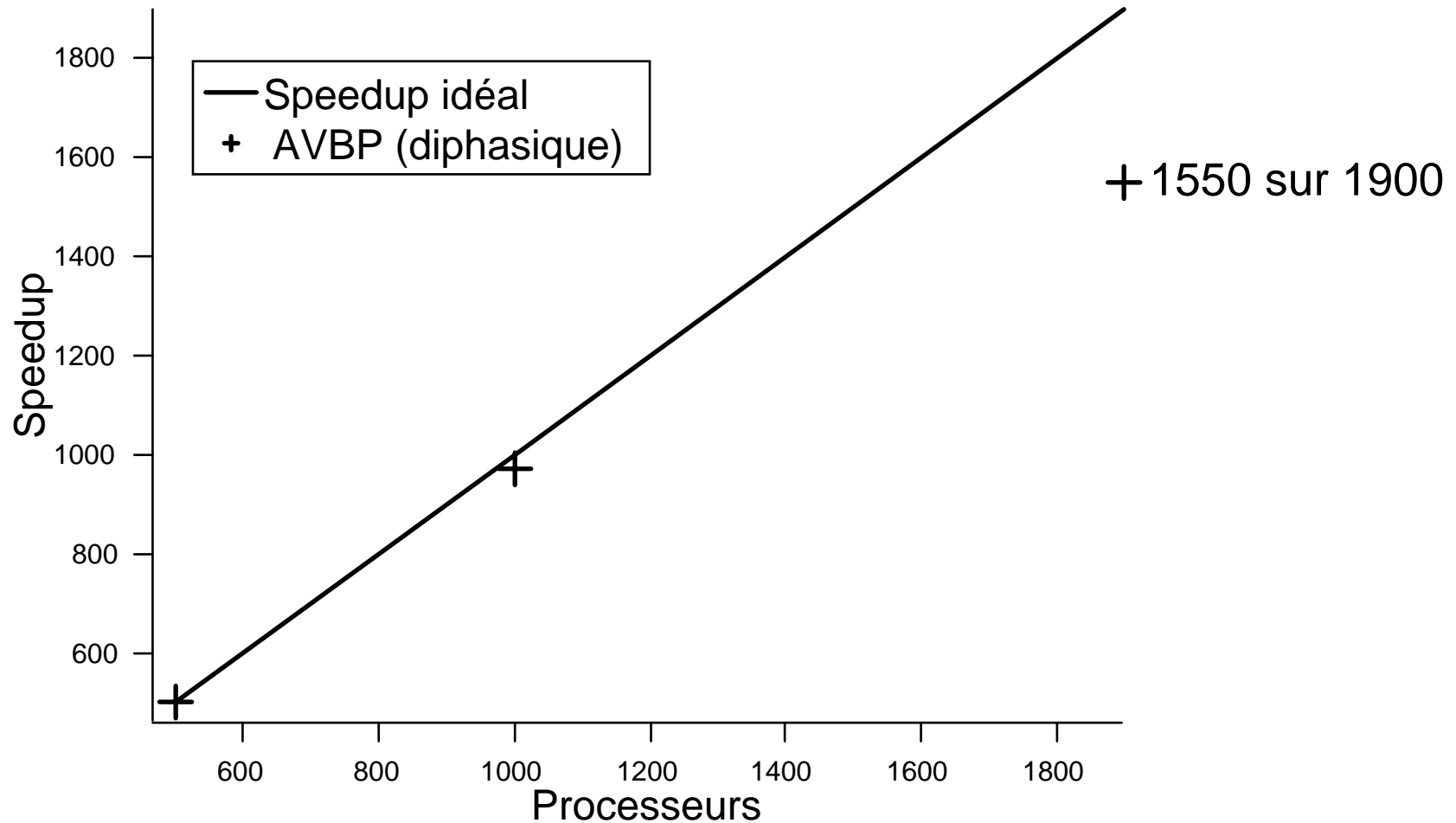
Accès à ter@10 1900 cores Montecito
BULL

L'évaporation est un phénomène très important si on souhaite évaluer les délais d'allumage

=> 2 fois plus d'équations à transporter!

QuickTime™ et un
décompresseur TIFF (LZW)
sont requis pour visionner cette image.

Pas encore de résultats à montrer (le calcul tourne actuellement),
juste un premier test de scalabilité:



Conclusion:

Ordre de grandeur des temps de restitution (extrapolation!) du calcul d'allumage gazeux Vesta

	Nombre de processeurs						
	1	32	120	500	1900	4096	5120
CRAY XD1	3 ans	1 mois	8 jours	-	-	-	-
BlueGene/L	11 ans	-	-	8 jours	2 jours	1 jour	19 heures
ter@10	2 ans 1/2	-	-	2 jours	11 heures	-	-
SGI 03800	11 ans	4 mois	1 mois	8 jours	-	-	-

Résumé:

|

	Code	nb procs	speedup
EDF	Saturne	1024	614
	Zephyr	128	104
	Dymoka	1024	1003
	VASP	32	5
CEA	TRIO U	512	282
CERFACS	AVBP	4096	4078