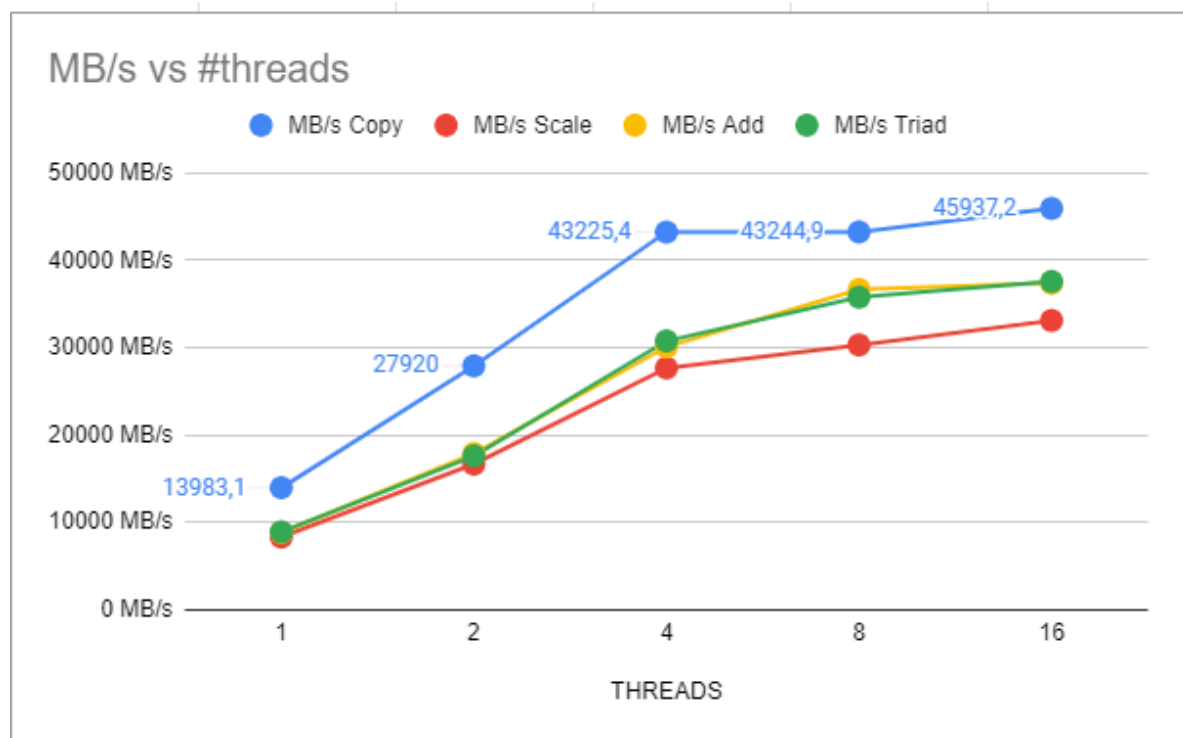


## LAB 3: Performance characterization of HPC clusters

En el document s'inclouen els resultats de diferents programes de prova. Els resultats es comparen amb la màquina BOADA. Cada programa de proves està dissenyat per un supòsit.

### Stream

Aquest joc de proves comprova operacions d'alt ús de memòria. L'ample de banda obtingut de les diferents operacions es representa en el gràfic.

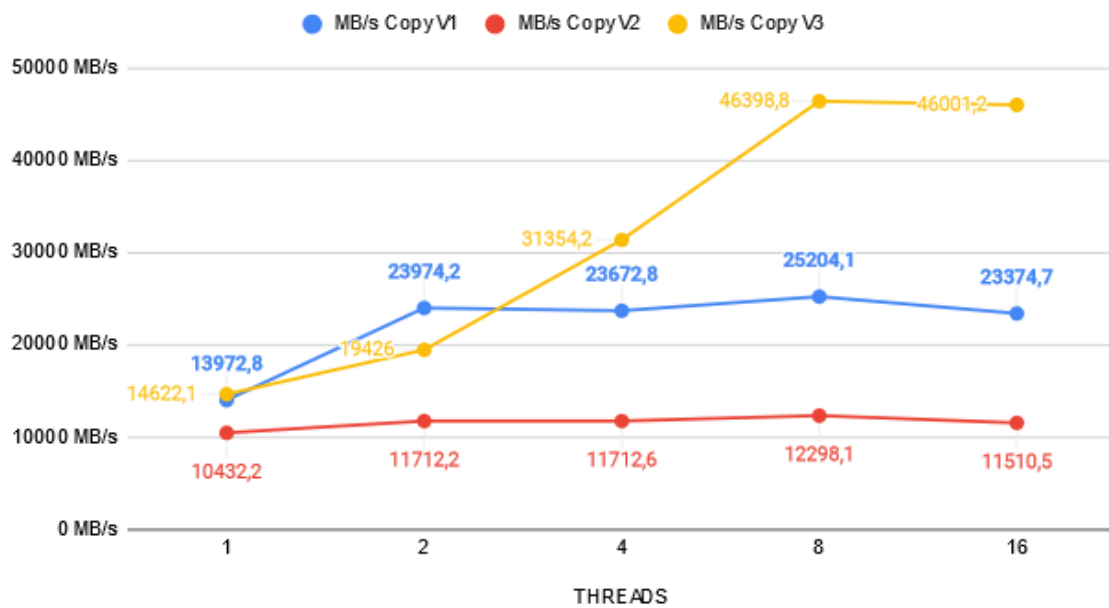


Pel que fa a l'ús de les directives numactl, tenim que aquestes controlen les polítiques de planificació i ubicació de memòria. L'opció `--membind` fa que nomès s'ubiqui memòria als nodes; d'altra banda que l'opció `--cpunodebind` fa que s'executin les comandes als nodes seleccionats.

S'han aplicat els següents paràmetres en els diferents test proposats:

- Ubicar la memòria i els threads a un NUMA node.(V1 al gràfic) `numactl -m 0 -N 0`
- Ubicar la memòria a un node i els threads a un altre node.(V2 al gràfic) `numactl -m 0 -N 1`
- Ubicar memòria i threads a dos nodes.(V3 al gràfic) `numactl -m 0,1 -N 0,1`

### Comparativa operació 'copy'



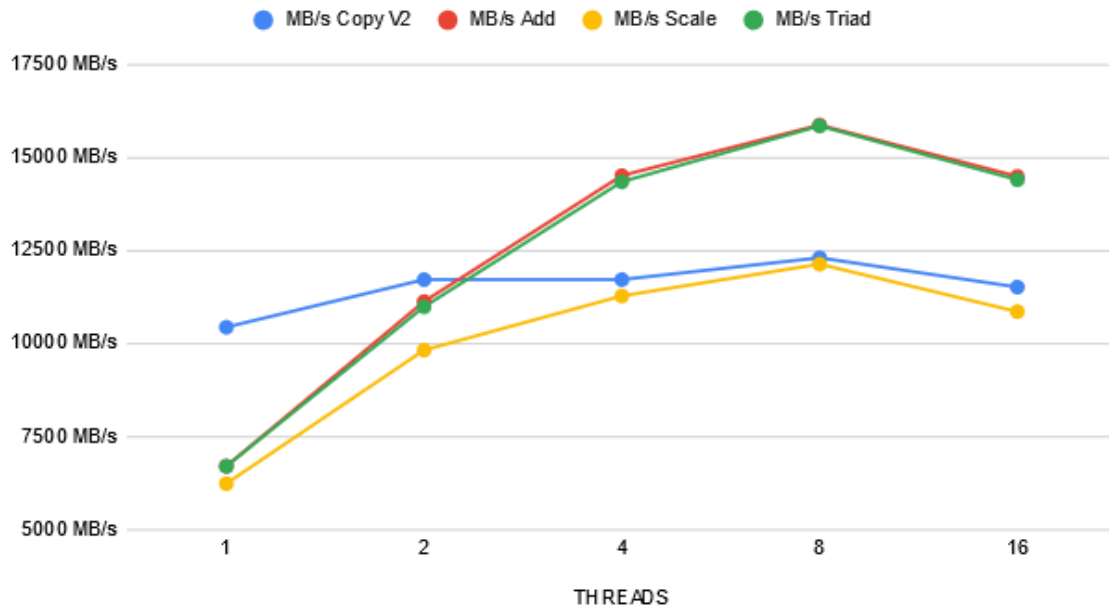
En el gràfic, que compara l'ample de banda en les tres versions podem veure:

- La segona versió (memoria i thread en diferent) node es la que pitjor rendiment obté, el rendiment ve condicionat per la comunicació de dades entre els threads i la memòria (que estan a diferent node).
- La primera versió no escala al augmentar el numero de threads, ja que no s'aprofita tot el paral·lelisme possible (respecte a la versió 3) al mantenir l'execució en un únic node.
- La tercera versió es la que dona el millor rendiment, molt similar al de la segona versió amb pocs threads pero ràpidament escala el seu rendiment amb els threads.
- Els resultats entre la primera i la tercera versió ens dona els resultats de casi el doble de rendiment (~1,84 i 1,96 vegades), fet esperat, ja que amb dos nodes tenim el doble de capacitat de memòria.

Per sota de 4 threads, no es compleix, segurament el sobrecost de la descomposició de tasques fa baixar el rendiment.

Respecte aquest segon gràfic, podem veure que les operacions que requereixen càlcul (add/triad) tenen un millor ample de banda que la operació de copia, això es dona pel temps de latència de l'operació, aquest temps emmascara la latència de memòria (que en aquest cas esta a un altre node). Es a dir, mentres es fa l'operació s'aprofita aquest temps per transmetre a l'altre node els resultats. Aquesta reflexió és interessant pero tot i això el rendiment d'aquesta versió es molt menor a les demés.

## Memòria i Threads a diferents nodes



Boada utilitza les memòries DDR4-2400 i DDR4-2133 que tenen sobre el paper un ample de banda de pic de 19.200 i 17.066 MB/s. En el millor dels casos (operació de copia) trobem que el rendiment màxim ronda els 46.000 MB/s.

## Linpack

Aquest *benchmark* està destinat a mesurar el rendiment en operacions de coma flotant. S'han mesurat el programa sobre 1 i 2 nodes amb diferents configuracions i nombre de processos. Als gràfics es veuen en verd les millors configuracions per cada configuració i l'*speed-up* respecte la pitjor configuració (també s'adjunta el full de càlcul).

- Per 1 node, s'han trobat els següents resultats:

1				2				4				8			
NB's	Time 1x1	GFlops	Speed-UP	NB's	Time 2x1	GFlops	Speed-UP	NB's	Time 4x1	GFlops	Speed-UP	NB's	Time 8x1	GFlops	Speed-UP
32	66.73	5.11	0	32	34.22	9.97	0	32	17.95	19.01	0	32	10.15	33.55	1.08
64	56.06	6.09	19.03	64	29.83	11.44	14.72	64	15.75	21.67	13.97	64	8.91	38.3	13.92
128	53.85	6.34	23.92	128	28.51	11.97	20.03	128	15.53	21.98	15.58	128	9.15	37.3	10.93
256	53.19	6.42	29.53	256	28.90	11.81	18.41	256	16.41	20.81	9.38	256	10.28	33.26	0.00

2				4				8				16			
NB's	Time 2x2	GFlops	Speed-UP	NB's	Time 4x2	GFlops	Scalability	NB's	Time 8x2	GFlops	Scalability	NB's	Time 16x2	GFlops	Scalability
32	17.31	19.72	3.70	32	9.24	36.95	11.04	32	5.07	42.32	27.14	32	2.77	40.53	21.85
64	15.19	22.47	18.17	64	8.14	41.97	26.04	64	2.81	42.00	26.20	64	1.51	39.49	18.75
128	14.85	22.95	20.88	128	8.19	41.71	25.27	128	1.51	38.54	15.80	128	0.84	34.98	5.02
256	15.44	22.11	16.26	256	8.87	38.50	15.97	256	0.88	38.54	15.80	256	0.77	34.98	5.02

- Per 2 nodes, s'han trobat els següents resultats:

2x1				1x2				4x1				2x2				1x4			
NB's	Time	GFlops	Speed-Up	Time	GFlops	Speed-Up		NB's	Time	GFlops	Speed-Up	Time	GFlops	Speed-Up		Time	GFlops	Speed-Up	
32	32.91	10.37	0.00	31.66	10.78	3.95		32	17.63	19.37	0	17.42	19.59	1.21		17.24	19.8	2.26	
64	29.6	11.53	11.18	28.68	11.90	14.75		64	15.69	21.76	12.36	15.19	22.47	16.06		15.18	22.49	16.14	
128	25.57	11.95	20.71	27.96	12.23	18.96		128	15.48	22.05	13.89	15.05	22.89	16.59		15.09	22.62	16.83	
256	28.89	11.81	13.91	28.25	12.07	16.37		256	16.37	20.85	7.70	15.45	22.09	14.11		15.98	21.36	10.33	

8x1				2x4				4x2				1x8			
NB's	Time	GFlops	Speed-Up	Time	GFlops	Speed-Up		Time	GFlops	Speed-Up		Time	GFlops	Speed-Up	
32	10.02	34.06	2.45	9.13	37.95	12.38		9.17	37.22	11.85		9.22	37.01	11.28	
64	8.93	38.22	14.89	8.10	42.16	26.67		8.68	42.25	26.96		8.38	40.73	22.43	
128	9.37	36.45	9.50	8.15	41.87	25.89		8.14	41.93	26.04		8.63	39.55	18.89	
256	10.26	33.26	0.00	8.85	38.57	15.93		8.85	38.57	15.93		9.75	35.00	5.23	

12x1				2x6				4x3				3x4				6x2				1x12			
NB's	Time	GFlops	Speed-Up	Time	GFlops	Speed-Up		Time	GFlops	Speed-Up		Time	GFlops	Speed-Up		Time	GFlops	Speed-Up		Time	GFlops	Speed-Up	
32	7.35	46.47	13.74	6.32	54.08	32.28		6.32	54.06	32.28		6.34	53.89	31.86		6.49	52.58	28.81		6.62	51.55	26.28	
64	6.70	50.88	24.78	5.73	59.57	45.90		5.62	59.25	46.15		5.64	60.54	46.23		5.82	58.05	43.64		6.13	55.86	36.38	
128	7.09	49.24	18.00	5.98	59.24	42.66		5.9	57.85	41.69		5.97	57.23	40.03		6.47	52.75	29.21		6.78	43.91	7.46	
256	8.36	40.84	0	6.54	52.20	27.83		6.38	52.25	31.03		6.54	52.93	27.83		6.69	51.03	24.96					

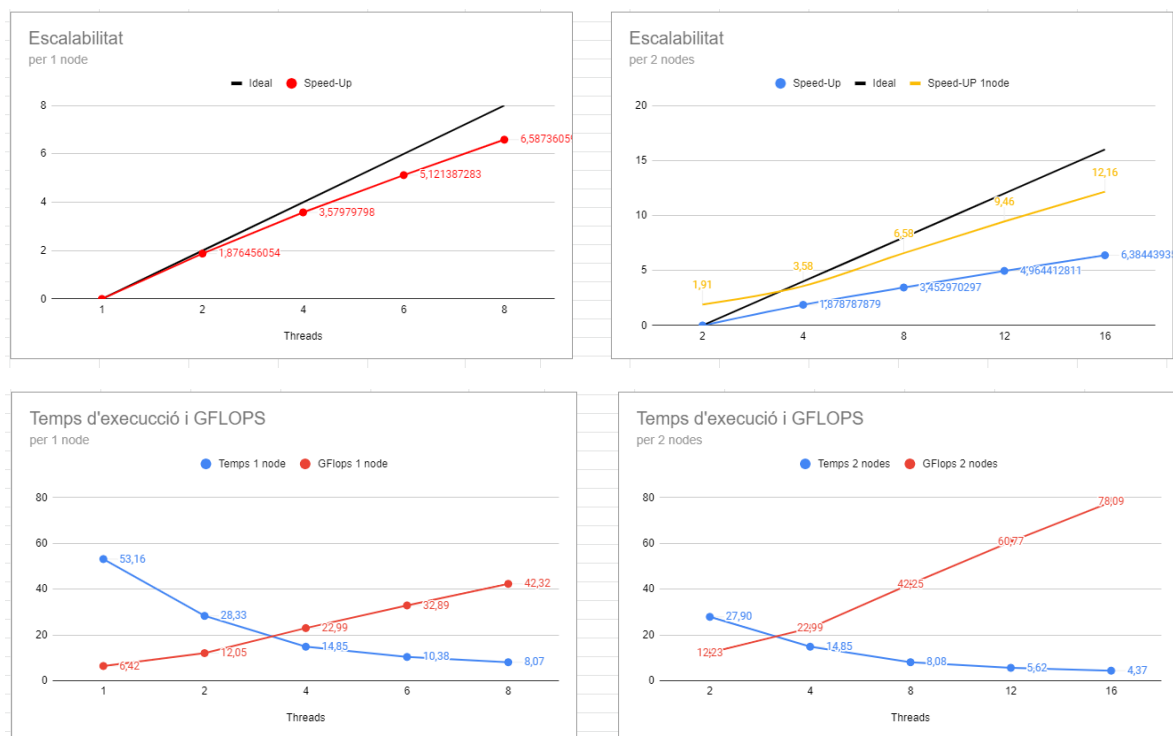
16x1				2x8				4x4				8x2				1x16			
NB's	Time	GFlops	Speed-Up	Time	GFlops	Speed-Up		Time	GFlops	Speed-Up		Time	GFlops	Speed-Up		Time	GFlops	Speed-Up	
32	6.23	54.76	18.39	4.93	69.24	49.49		4.94	69.11	49.19		5.25	65.06	40.38		5.32	64.23	38.53	
64	5.66	60.27	30.21	4.53	75.34	62.69		4.37	76.89	66.95		4.68	72.93	57.48		5.02	68.00	46.81	
128	6.19	55.14	19.06	4.71	72.49	56.48		4.51	75.63	63.41		4.88	69.92	51.02		5.41	63.05	36.23	
256	7.37	46.35	0	5.45	62.70	35.23		5.12	66.74	43.95		5.68	60.12	29.75		6.57	51.95	12.18	

Podem comprovar que les configuracions que millor funcionen són aquells on la graella de processos està "equilibrada", ja que així la càrrega de treball es distribueix millor i obtenim més rendiment.

També podem veure que la mida òptima està entre 64 i 128, podem fer les següents hipòtesis:

- En el cas d'un procés i 1 node, obtenim més rendiment si fem menys divisions.
- Per sota de 8 processos, la mida òptima és 128, degut possiblement a què cada procés pot reservar aquest espai de memòria sense tenir conflictes per l'espai.
- Amb 8 processos la mida baixa a 64, seguint amb el raonament anterior, es possible que els diferents processos en reservar el mida de 128 tinguin conflictes i es perdi rendiment per culpa de la distribució de la memòria.

Agafant els millors resultats obtenim els següents gràfics per un i dos nodes on podem veure l'escalabilitat i el rendiment:



Cada node Boada té 2 [Intel Xeon E5-2609 v4 @ 1,70GHz](#), amb 8 unitats de càlcul a 1,70 GHz i instruccions multimedia (Intel AVX2 capaç de fer operacions de 256 bits i FMA). Per tant calculem el pic de GFLOPS/s de cada node com:

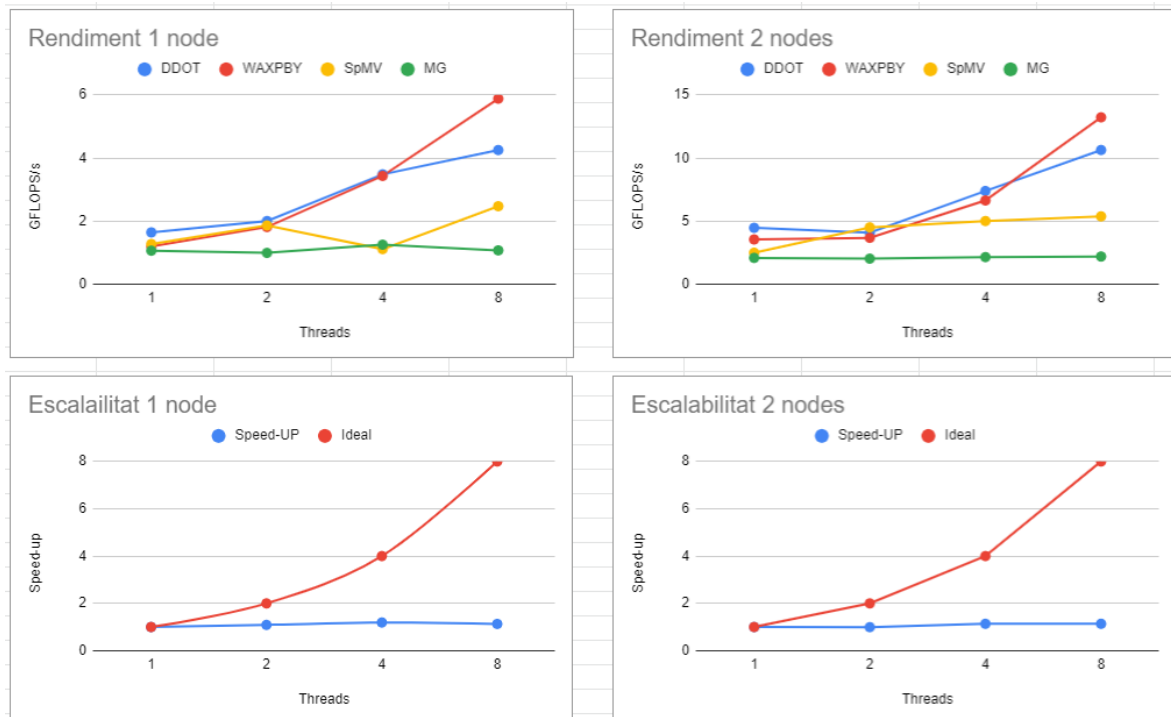
$$\text{GFLOPS de pic} = 8 \text{ unitats} * (256/64) * 2 \text{ flops} * 1,70 \text{ Gcycles\_per\_segons} \\ = 1.088 * 10^{11} = 108.8 \text{ GFLOPS}$$

Mirant els resultats, podem veure que el nostre rendiment en GFLOPS/s es de 42,32 (per un node) i 78,09 (per dos nodes). S'ha aconseguit un 38,9% i 35,88% del pic respectivament. Aquest resultats estan afectats pels temps/latències de memòria i comunicació entre nodes. Es pot concloure que l'intensitat aritmètica per cada *byte* es bastant petita.

## HPCG

Aquest últim programa està orientat a mesurar el rendiment sobre aplicacions irregulars. Està format per 4 *kernels* que fan operacions diferents

Prenent les mesures de rendiment, obtenim les següents gràfiques. L'escalabilitat està calculada sobre el



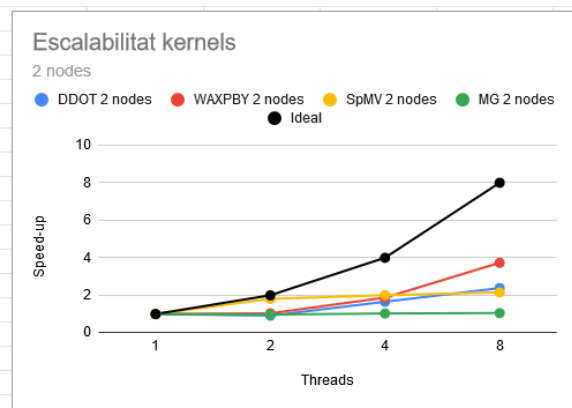
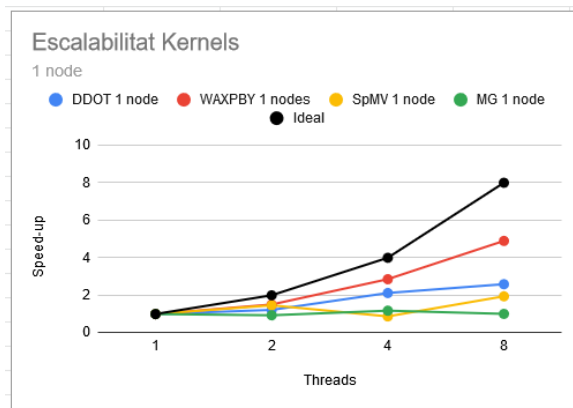
Podem veure que l'escalabilitat total; com era d'esperar, queda lluny de l'ideal i el rendiment es bastant irregular.

El càlcul de GFLOPS/s total, es fa ponderant el temps d'execució de cada *kernel*. Aproximadament tenim la següent ponderació per cada *kernel*:

- DDOT(0,008%): Producte de punts de dos vectors.
- WAXPBY(0,009%): Actualització d'un vector amb la suma de dos vectors.
- SpMV(0,09%): Multiplicació de matrius esparses
- MG(0,89%): Simètric Gauss-Seidel

Veient això, si haugéssim de millorar el rendiment, hauríem de millor aquella secció de codi més utilitzada (lleï d'Amdahl).

També podem observar que l'escalada de cada kernel, exceptuan WAXPBY i el DDOT, tenen un rendiment força irregular. També podem veure, en general que queda lluny de l'escalat ideal.



Comparant els resultats d'aquest *benchmark* amb el Linpack(agafant el temps d'execució), obtenim:

- Rendiment Linpack 2 nodes i 8 Threads:  $42,25 \text{ GFLOPS} / 8,08 \text{ s} = 5,22 \text{ GFLOPS/s}$
- Rendiment HPCG 2 nodes i 8 Threads :  $2,33 \text{ GFLOPS/s}$

Podem veure com el rendiment del Linpack es superior, era d'esperar ja que el Linpack es un *benchmark* on sempre es fa la mateixa operació, en contra del HPCG.