

1 Descrizione dell'applicazione

È stata realizzata una applicazione sul calcolo del prodotto matrice per vettore. La computazione considerata è su stream, gli elementi dello stream sono le matrici su cui effettuare il calcolo, mentre il vettore rimane costante per tutta l'esecuzione dell'applicazione. L'implementazione scelta fa uso del paradigma Data Parallel, in particolare la forma *Map*, in quanto ogni processo worker andrà ad eseguire un calcolo che è indipendente dall'attività degli altri processi. Si è scelto il partizionamento delle matrici per riga. Per realizzare le comunicazioni vengono usati i canali del supporto fornito, quindi si ha lo scambio dei puntatori alle strutture dati condivise. L'applicazione risulta seguire lo schema *multicast-compute-gather* in quanto per ogni matrice dello stream si eseguono le seguenti tre fasi:

- a) la distribuzione del riferimento alla matrice corrente tra i processi worker, è compito di ogni worker effettuare il calcolo nella propria partizione della matrice,
- b) il calcolo dei risultati parziali in ciascun worker,
- c) la raccolta di tutti i risultati parziali dell'elemento.

La comunicazione collettiva *multicast* è implementata con una struttura ad albero binario mappato nell'insieme dei processi worker, le comunicazioni che costituiscono l'albero sono implementate per mezzo dei canali simmetrici offerti dal supporto; l'altra comunicazione collettiva, la *gather*, viene implementata per mezzo del canale asimmetrico in ingresso fornito dal supporto.

L'applicazione è fittizia, nel senso che non esistono dispositivi che generino e collezionino lo stream, per questo motivo l'applicazione è costituita, oltre che dai processi che realizzano la *map*, da altri due processi che, rispettivamente, generano e collezionano gli elementi dello stream. Tali due processi comunicano con il sottosistema dei workers ($\Sigma 1^{(n)}$) per mezzo dei canali del supporto: il processo generatore è collegato tramite un canale simmetrico al worker radice dell'albero multicast, ogni processo worker è collegato al processo collettore per mezzo di un canale asimmetrico.

Riassumendo, ogni processo worker fa uso di 4 canali di comunicazione:

- tre sono simmetrici e trasportano gli elementi dello stream, uno dei quali è in ingresso dal processo padre nell'albero della multicast, gli altri due sono in uscita verso i processi radice dei due sottonodi della multicast,
- il quarto canale è asimmetrico in ingresso ed è usato in scrittura, per la comunicazione del risultato parziale al processo collezionatore.

È quindi possibile l'uso del supporto alle comunicazioni che fa uso della UDN in quanto, per ogni processo, il numero di canali non supera il numero di code hardware.

2 Analisi delle prestazioni

Di seguito viene descritta una breve analisi delle prestazioni attese dal benchmark. Dato che la frequenza dello stream è arbitraria, la stima del tempo di servizio del *map* permette di avere un primo dimensionamento del tempo di interarrivo, per diverse dimensioni dei dati, in modo tale da poter effettuare le prime misure dell'applicazione. Nella sezione successiva sono quindi proposte le misure del tempo di completamento dello stream e del tempo di servizio di $\Sigma 1^{(n)}$.

La macchina *TILEPro64* non dispone di processori di comunicazione, ne segue che le latenze di comunicazione dei canali sono pagate completamente nel tempo di servizio dei processi worker del sottosistema $\Sigma 1^{(n)}$. Si caratterizza perciò il tempo di servizio *ideale* ed *effettivo* del sottosistema *map* come segue:

$$T_{\Sigma 1_id}^{(n)} = T_S^{(n)} = T_{multicast} + \frac{T_{cagl}}{n} + T_{gather} = 2 \cdot T_{sym_send} + \frac{T_{cagl}}{n} + T_{asym_send} = \Delta + \frac{T_{cagl}}{n}$$

$$T_{\Sigma 1}^{(n)} = \max(\{T_A, T_S^{(n)}\})$$

Dove T_{calc} è il tempo medio impiegato per il calcolo della computazione sequenziale, ovvero il calcolo di una moltiplicazione matrice per vettore, e n è il grado di parallelismo dell'applicazione, inteso come il numero dei processi worker. Il rapporto tra T_{calc} e n esprime il tempo di servizio *ideale* di un processo worker scollegato dallo stream. Si è indicato con Δ la somma dei tempi spesi nelle comunicazioni da parte di un worker, tale latenza non è sovrapposta al tempo di calcolo nei processi worker.

<i>Matrix Size</i>	$T_{\text{calc}} (\mu s)$	$T_{\text{calc}} (\tau)$
56x56	85.997340	74351.900000
168x128	848.096504	733250.424000
280x280	2360.060404	2040469.784000

Table 1: Tempi di calcolo della computazione sequenziale al variare della dimensione della matrice

Se è noto il valore del T_{calc} e di Δ allora per un generico tempo di interarrivo si ricava il grado di parallelismo ottimo, ovvero il minor grado di parallelismo che massimizza l'efficienza fornendo un fattore di utilizzazione unitario del sottosistema $\Sigma 1^{(n)}$:

$$n_{\text{opt}} = \min(\{n \in \mathbb{N} \mid T_S^{(n)} \leq T_A\}) = \left\lceil \frac{T_{\text{calc}}}{T_A - \Delta} \right\rceil$$

Dato che il tempo di interarrivo non è definito a priori ma è arbitrario, da tale formula è possibile ricavare il tempo di interarrivo uguale al tempo di servizio del *map* con il massimo grado di parallelismo esplicitabile dall'architettura.

Si pone quindi il problema di stimare il tempo di calcolo e il valore di Δ .

- Il T_{calc} viene stimato misurando il tempo medio impiegato per eseguire una moltiplicazione matrice per vettore in un singolo tile della macchina. I risultati di tale misura per dimensioni diverse della matrice sono mostrate in Tabella ??.

- Il valore di Δ può essere fornito in prima approssimazione dalle misure delle latenze di comunicazione effettuate con l'applicazione “ping-pong” per le due implementazioni dei canali. I risultati di tale misura sono riassunti nella Tabella ??.

Si osserva che le misure di tali parametri sono approssimazioni ottimistiche, è infatti prevedibile che durante l'esecuzione dell'applicazione *map* sia il tempo di calcolo dei worker che le latenze di comunicazione siano superiori ai valori stimati per mezzo delle applicazioni di misurazione, le quali usano un numero minimale di processi. Per la misura del T_{calc} viene eseguito un unico processo, per la misura delle latenze di comunicazione vengono eseguiti due processi che si scambiano messaggi usando il supporto fornito, in entrambi i casi le misure sono prese con un numero di conflitti minimo sia nelle reti di interconnessione sia alle memorie cache e ai controllori della memoria principale. Durante l'esecuzione della *map* i processi in gioco possono essere molti fino al massimo numero di processori utilizzabili nella macchina, ne segue un aumento dei conflitti alle reti e alle memorie rispetto a quelli che si verificano nei programmi di misurazione e ciò introduce overheads sia nel tempo di calcolo effettivo dei workers che nelle latenze di comunicazione.

Si calcola il tempo di servizio ideale con il grado di parallelismo massimo, $N = 59$, $T_S^{(n)} = \frac{T_{\text{calc}}}{n} + \Delta$.

$$\begin{aligned} 74352/59 + 181 &= 1441 \quad \rightarrow 4000 \\ 74352/59 + 481 &= 1741 \quad \rightarrow 7000 \\ 74352/59 + 725 &= 1985 \end{aligned}$$

$$\begin{aligned} 733250/59 + 181 &= 12428 + 181 = 12609 \quad \rightarrow 20000 \\ &12428 + 481 = 12909 \quad \rightarrow 20000 \end{aligned}$$

$$2040470/59 + 181 = 34584 + 181 = 34765 \quad \rightarrow 50000$$

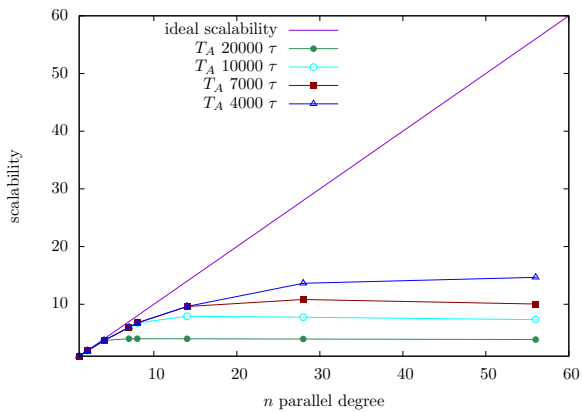
<i>Canale</i>	$L_{\text{com}} (\tau)$
ch_sym_udn	55.722760
ch_sym_sm_rdyack	155.377520
ch_asym_min_udn	69.509710
ch_asym_min_sm	170.285120
ch_asym_min_sm_all	414.157690

<i>Canali usati</i>	$T_{\Delta}(\tau)$
UDN only	180.95523
SM only	481.04016
SM only with all	724.91273

Table 2: Misure delle latenze dei canali di comunicazione rilevate con l'applicazione “ping-pong”, nella quale i due processi sono eseguiti in due tile con distanza massima nella mesh

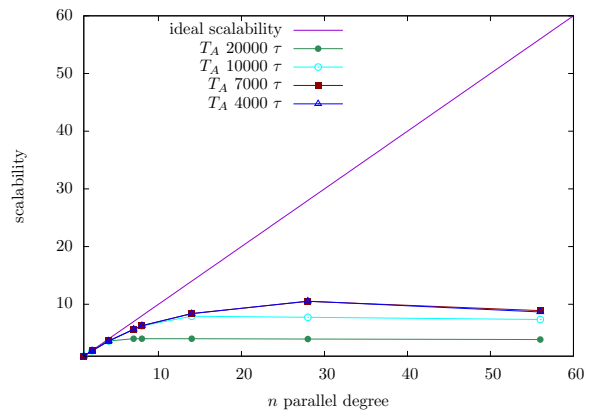
Figure 1: Grafici di scalabilità del tempo di completamento dello stream al variare del tempo di interarrivo

(a) Implementazione con solo UDN

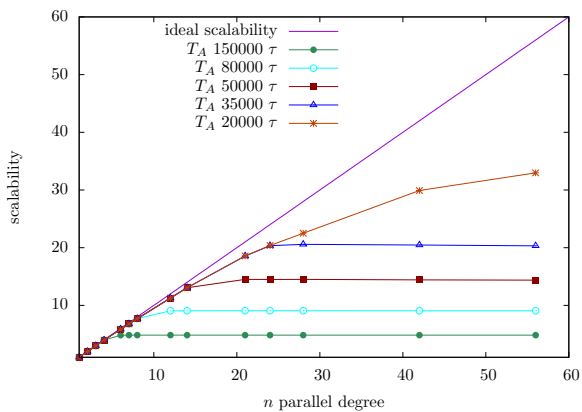


(a1) Scalabilità dell'implementazione UDN con M=56 al variare del tempo di interarrivo

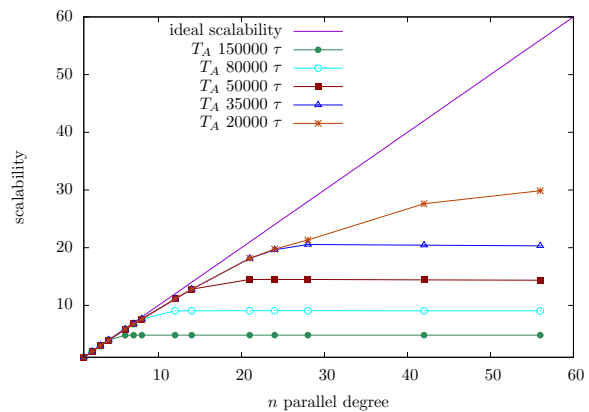
(b) Implementazione con solo SM



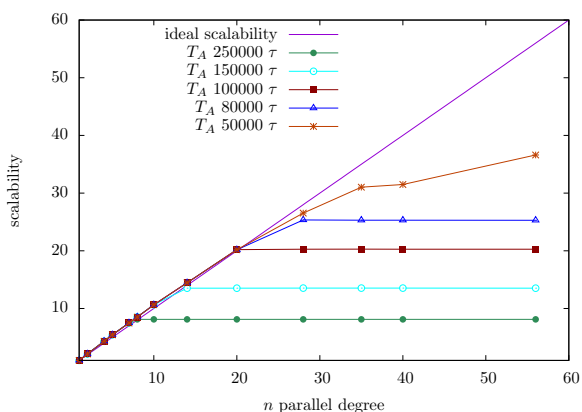
(b1) Scalabilità dell'implementazione SM con M=56 al variare del tempo di interarrivo



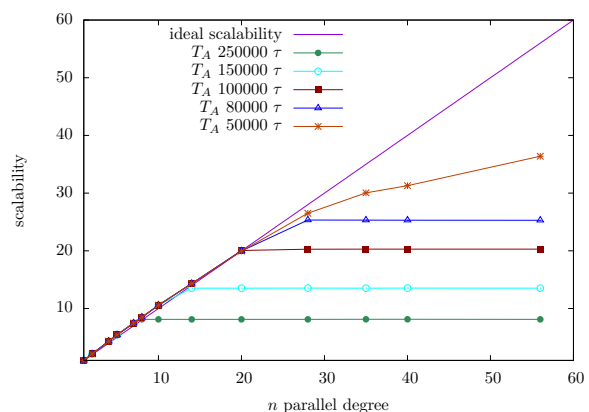
(a2) Scalabilità dell'implementazione UDN con M=168 al variare del tempo di interarrivo



(b2) Scalabilità dell'implementazione SM con M=168 al variare del tempo di interarrivo



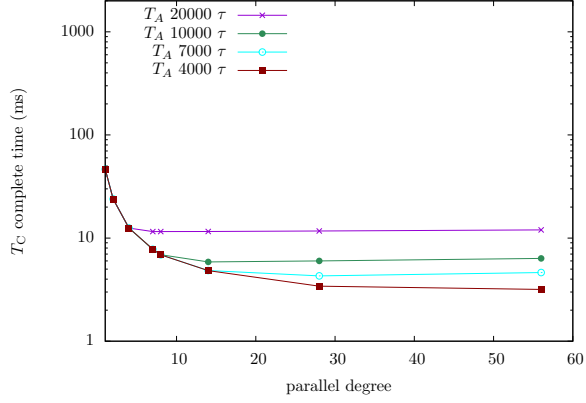
(a3) Scalabilità dell'implementazione UDN con M=280 al variare del tempo di interarrivo



(b3) Scalabilità dell'implementazione SM con M=280 al variare del tempo di interarrivo

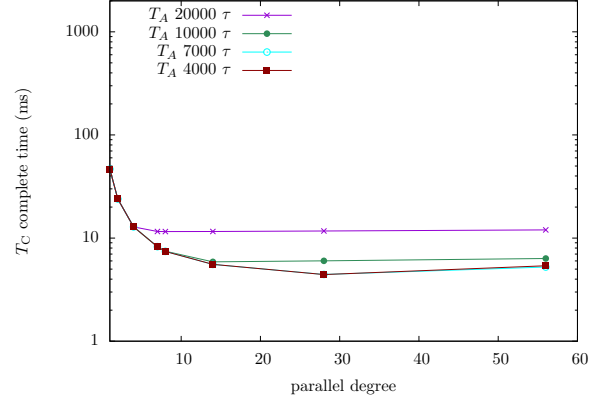
Figure 2: Grafici del tempo di completamento al variare del tempo di interarrivo

(a) Implementazione con solo UDN

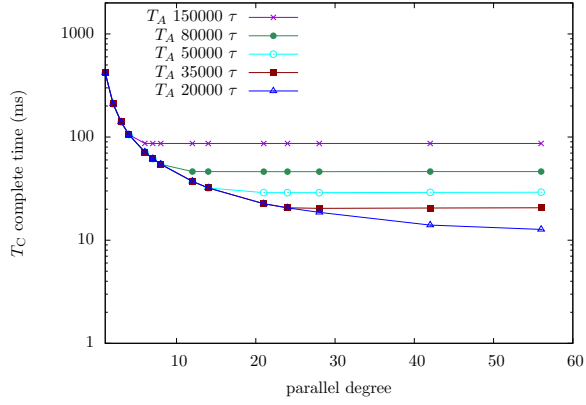


(a1) Tempo di completamento dell'implementazione UDN con $M=56$ al variare del tempo di interarrivo

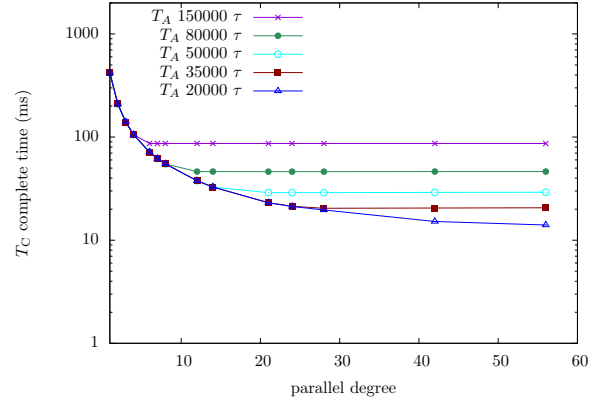
(b) Implementazione con solo SM



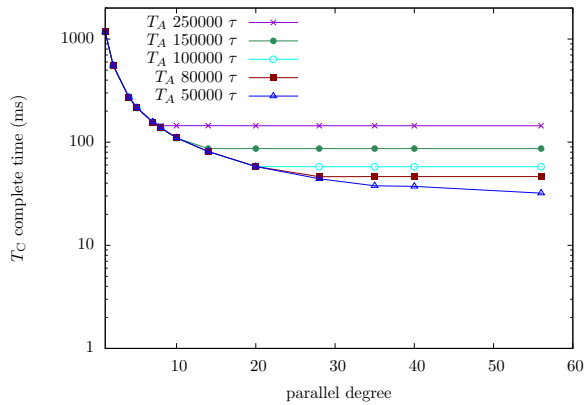
(b1) Tempo di completamento dell'implementazione SM con $M=56$ al variare del tempo di interarrivo



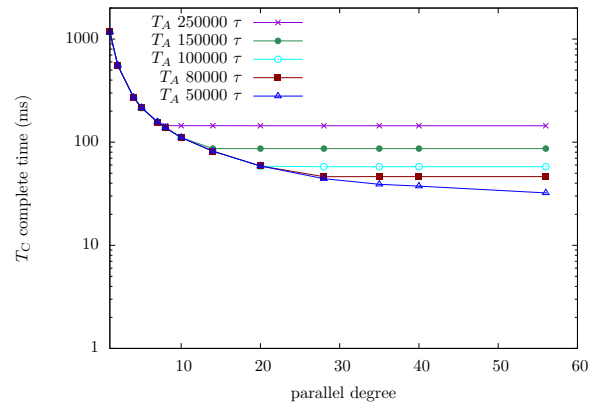
(a2) Tempo di completamento dell'implementazione UDN con $M=168$ al variare del tempo di interarrivo



(b2) Tempo di completamento dell'implementazione SM con $M=168$ al variare del tempo di interarrivo



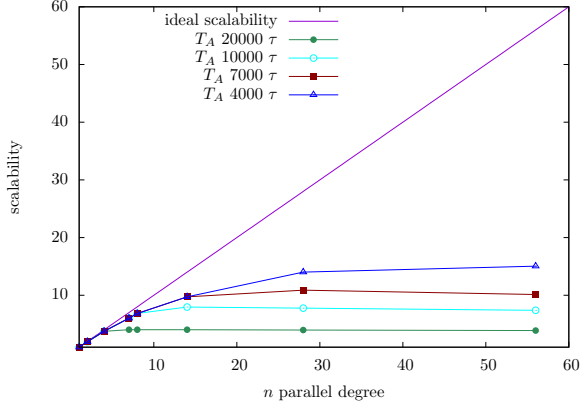
(a3) Tempo di completamento dell'implementazione UDN con $M=280$ al variare del tempo di interarrivo



(b3) Tempo di completamento dell'implementazione SM con $M=280$ al variare del tempo di interarrivo

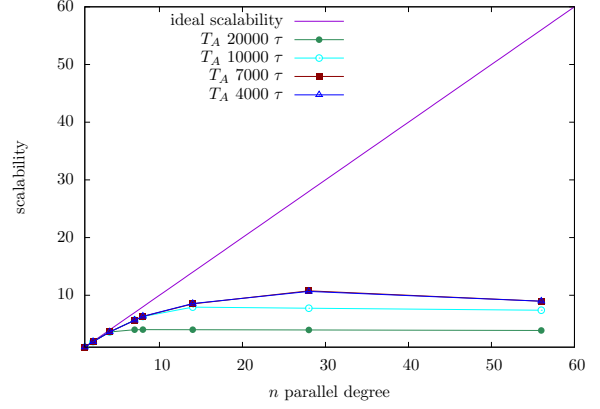
Figure 3: Grafici di scalabilità del tempo di servizio dello stream al variare del tempo di interarrivo

(a) Implementazione con solo UDN

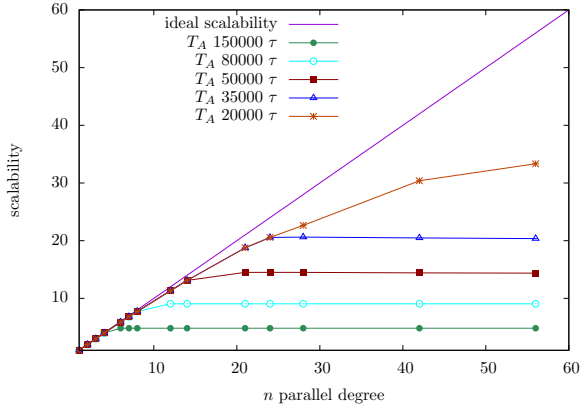


(a1) Scalabilità dell'implementazione UDN con M=56 al variare del tempo di interarrivo

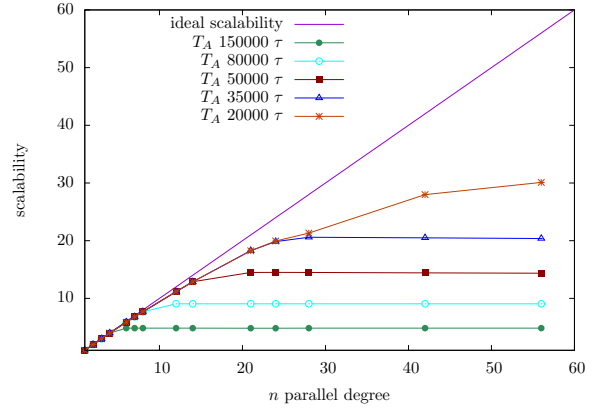
(b) Implementazione con solo SM



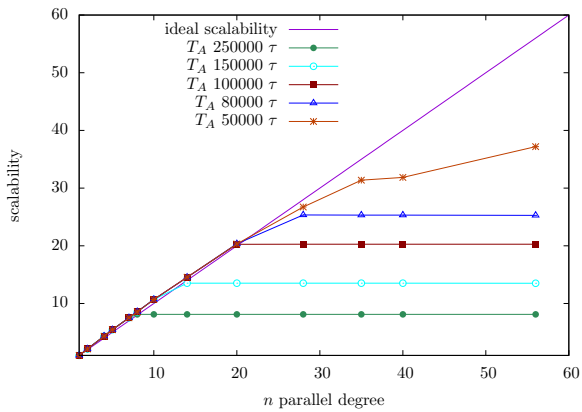
(b1) Scalabilità dell'implementazione SM con M=56 al variare del tempo di interarrivo



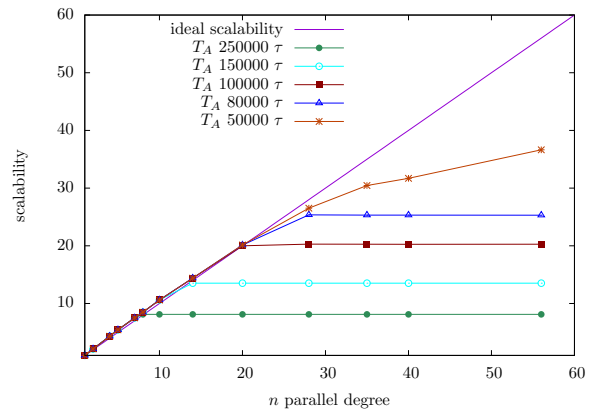
(a2) Scalabilità dell'implementazione UDN con M=168 al variare del tempo di interarrivo



(b2) Scalabilità dell'implementazione SM con M=168 al variare del tempo di interarrivo



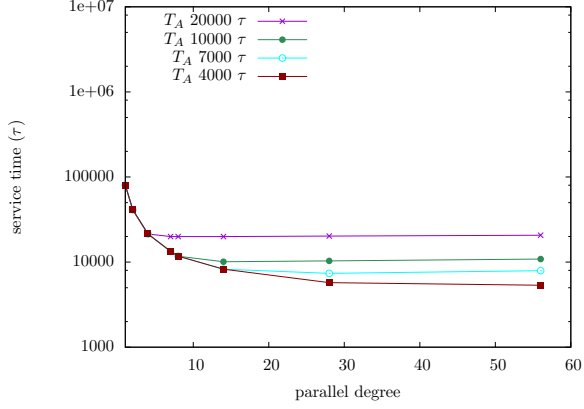
(a3) Scalabilità dell'implementazione UDN con M=280 al variare del tempo di interarrivo



(b3) Scalabilità dell'implementazione SM con M=280 al variare del tempo di interarrivo

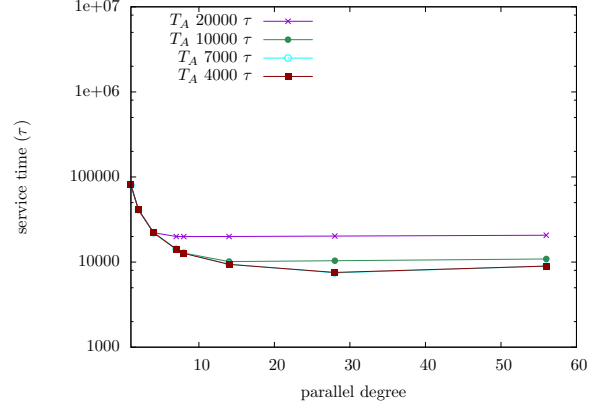
Figure 4: Grafici del tempo di servizio al variare del tempo di interarrivo

(a) Implementazione con solo UDN

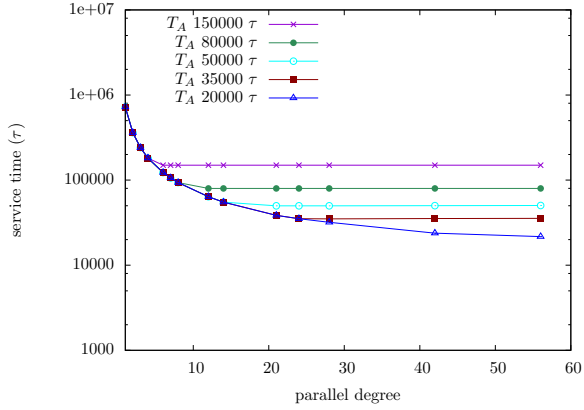


(a1) Tempo di servizio dell'implementazione UDN con M=56 al variare del tempo di interarrivo

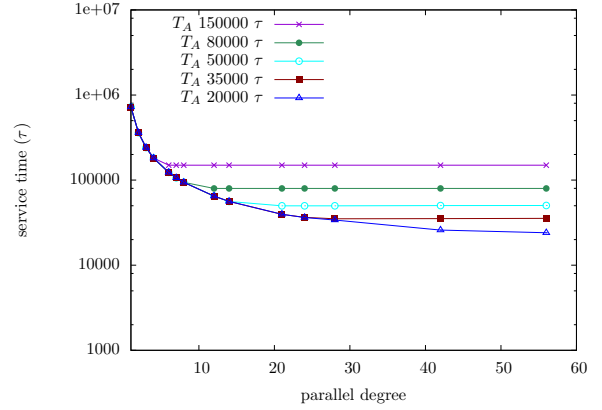
(b) Implementazione con solo SM



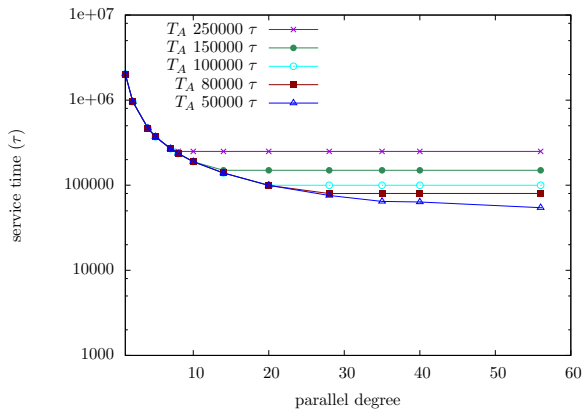
(b1) Tempo di servizio dell'implementazione SM con M=56 al variare del tempo di interarrivo



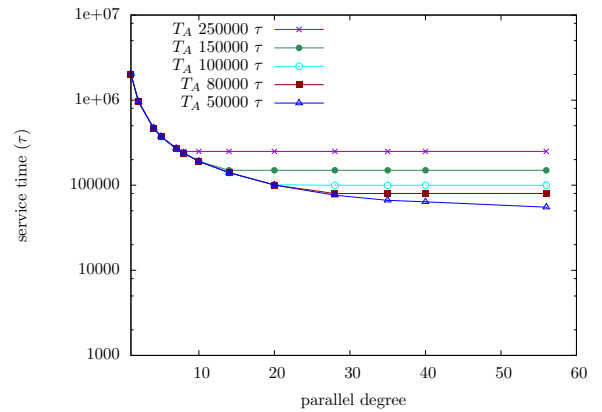
(a2) Tempo di servizio dell'implementazione UDN con M=168 al variare del tempo di interarrivo



(b2) Tempo di servizio dell'implementazione SM con M=168 al variare del tempo di interarrivo



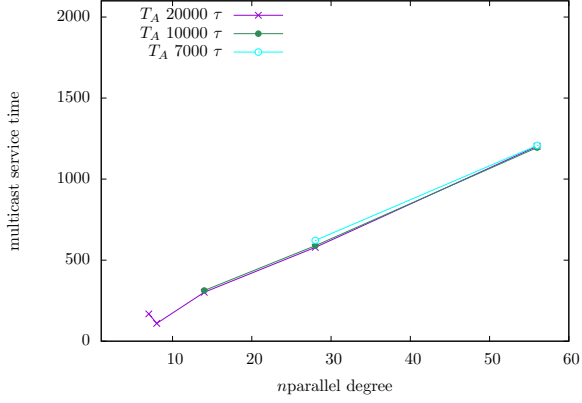
(a3) Tempo di servizio dell'implementazione UDN con M=280 al variare del tempo di interarrivo



(b3) Tempo di servizio dell'implementazione SM con M=280 al variare del tempo di interarrivo

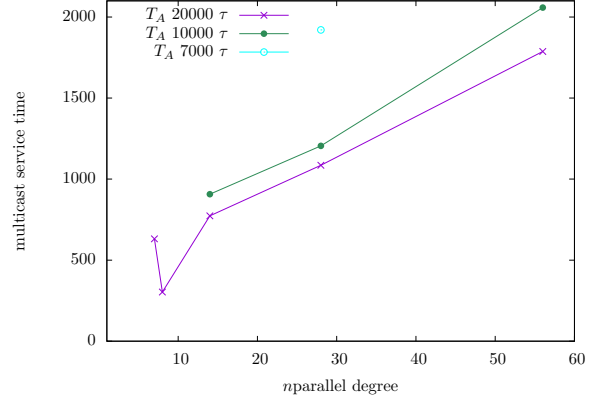
Figure 5: Grafici del tempo di multicast al variare del tempo di interarrivo

(a) Implementazione con solo UDN

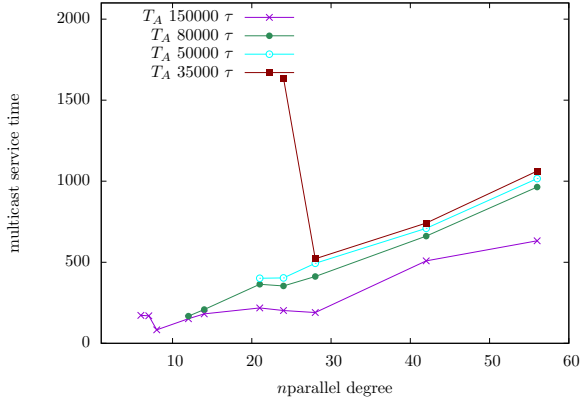


(a1) tempo di multicast dell implementazione UDN con M=56 al variare del tempo di interarrivo

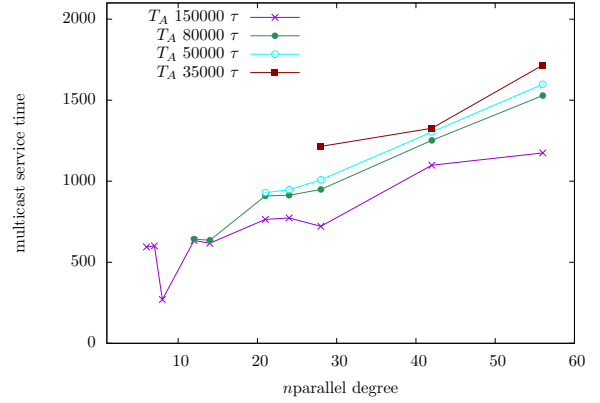
(b) Implementazione con solo SM



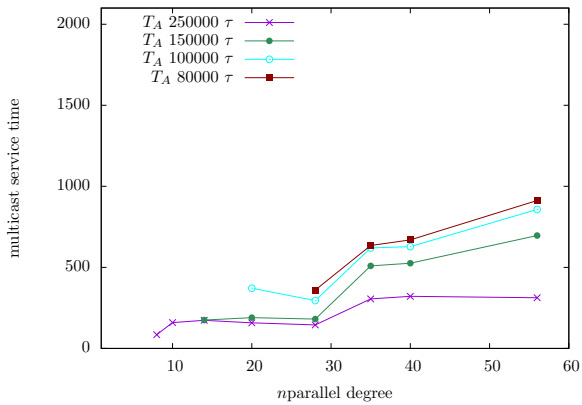
(b1) tempo di multicast dell implementazione SM con M=56 al variare del tempo di interarrivo



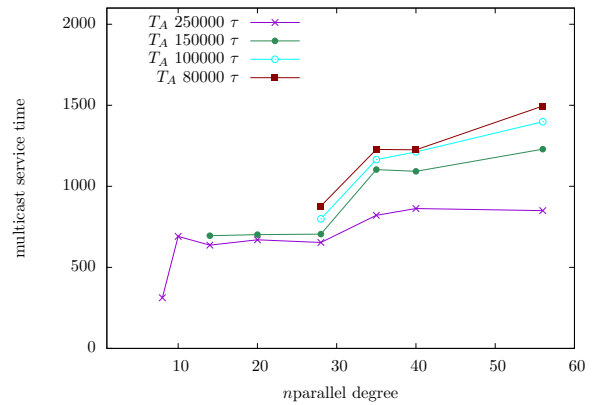
(a2) tempo di multicast dell implementazione UDN con M=168 al variare del tempo di interarrivo



(b2) tempo di multicast dell implementazione SM con M=168 al variare del tempo di interarrivo

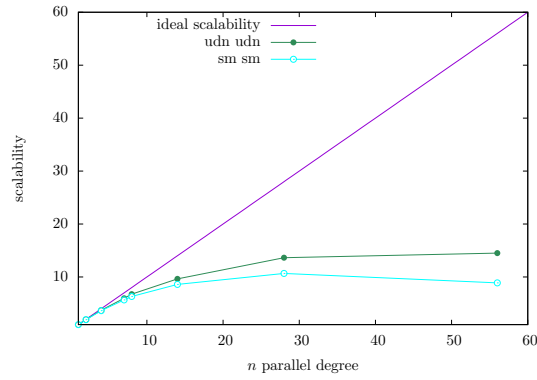


(a3) tempo di multicast dell implementazione UDN con M=280 al variare del tempo di interarrivo

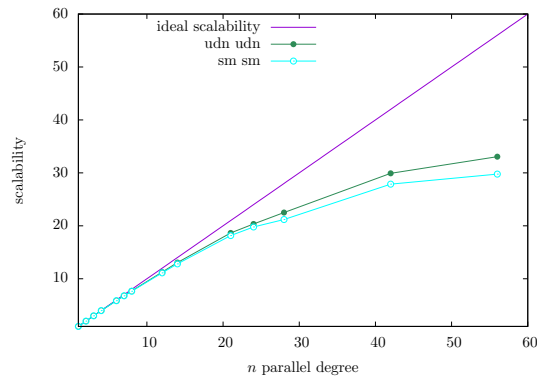


(b3) tempo di multicast dell implementazione SM con M=280 al variare del tempo di interarrivo

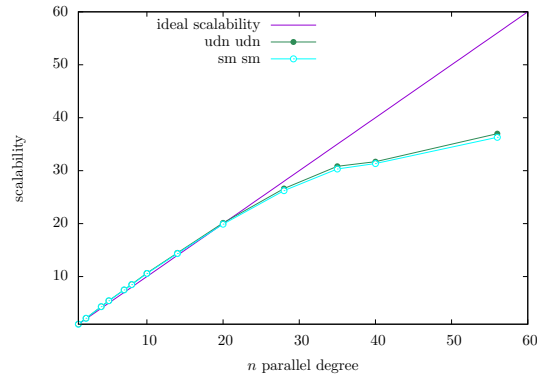
2.1 Confronto scalabilità delle due implementazioni



(a) Confronto della scalabilità nelle diverse implementazioni, $Ta=181$, $M=56$



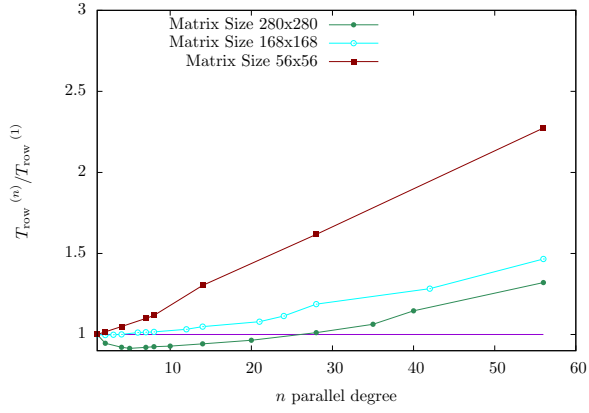
(b) Confronto della scalabilità nelle diverse implementazioni, $Ta=181$, $M=168$



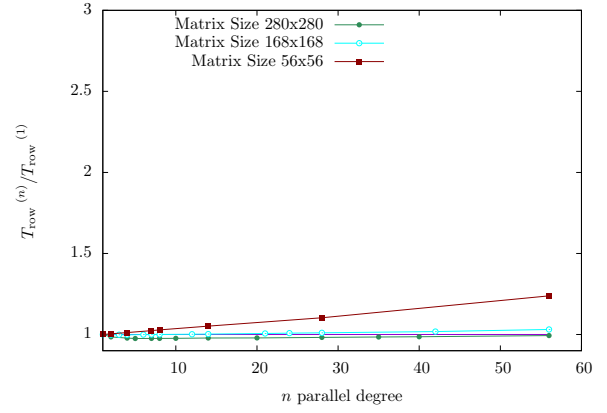
(c) Confronto della scalabilità nelle diverse implementazioni, $Ta=181$, $M=280$

2.2 Row calculation time

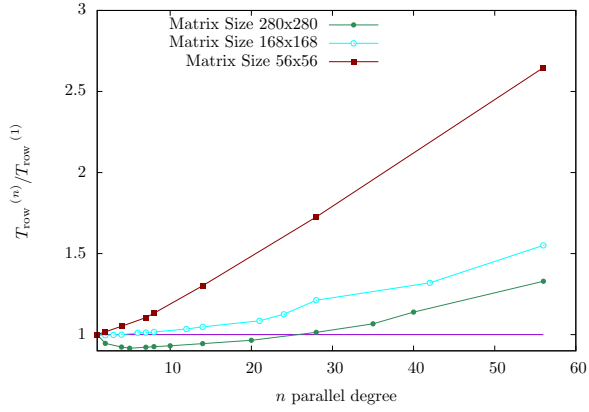
Figure 6: Rapporto tra i tempi di calcolo di un singolo prodotto scalare



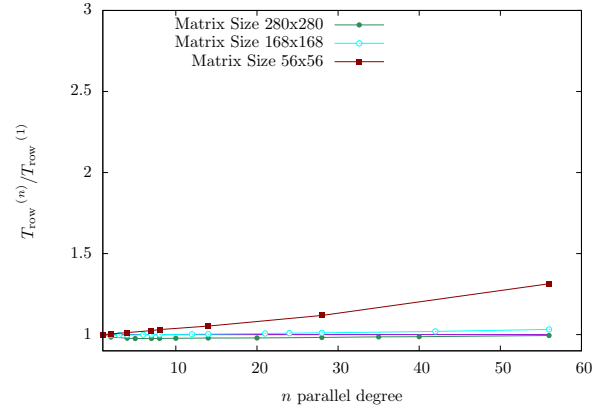
(d) Tempi di calcolo di una singola Row · Col con $T_a=4000$, canali UDN e dati Int



(e) Tempi di calcolo di una singola Row · Col con $T_a=4000$, canali UDN e dati Float



(f) Tempi di calcolo di una singola Row · Col con $T_a=4000$, canali SM e dati Int



(g) Tempi di calcolo di una singola Row · Col con $T_a=4000$, canali SM e dati Float