

Word Count with MPI

An MPI implementation of Word Count problem

Luigi Crisci

Carmine Spagnuolo

Abstract

Word Count rappresenta uno dei problemi più famosi nel calcolo distribuito, in cui tutti gli studiosi di tale settore si sono imbattuti almeno una volta nella loro vita. Ne verrà presentata un'implementazione in MPI, correlata da un'attenta analisi delle prestazioni.

1. Introduzione

Word Count consiste nel determinare in un testo quante occorrenze ci sono di ogni parola presente. Generalmente si presenta in applicazioni in cui è necessario mantenere le parole inserite in limiti definiti (es. nella narrativa, il numero di parole definisce la categoria dello scritto).

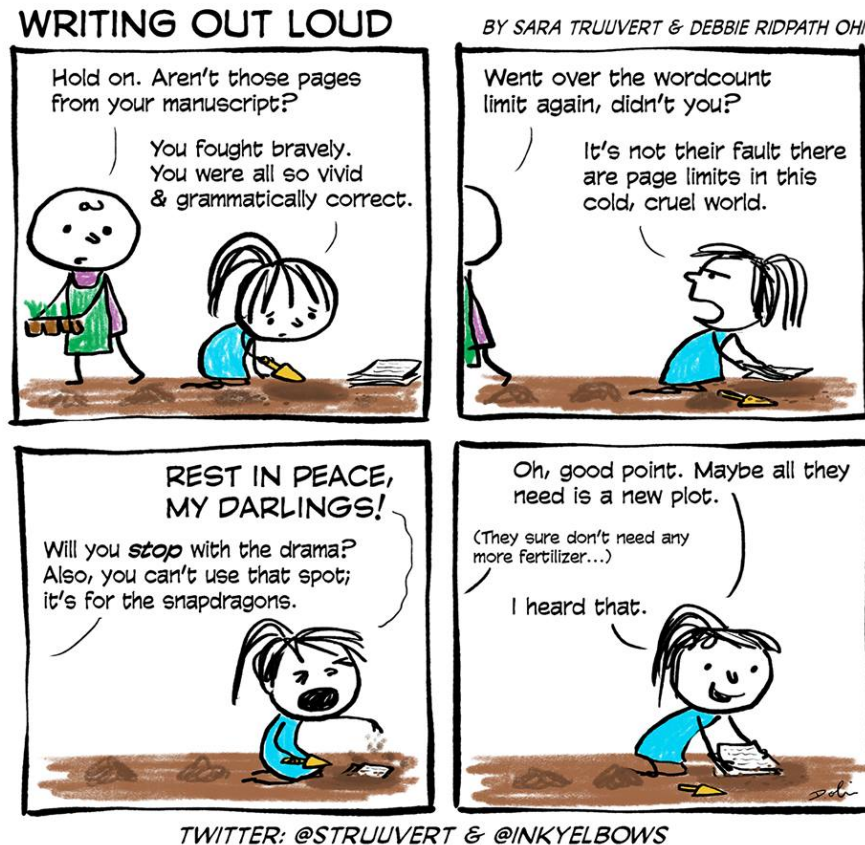


Figure 1: Word Count in letteratura

Benchè estremamente semplice nella sua definizione, Word Count rappresenta una sfida nel campo della programmazione distribuita a causa dell'enorme taglia dei problemi che possono presentarsi: pensiamo infatti a situazioni in cui i documenti da esaminare hanno dimensioni arrivanti anche a 100 TB, in cui un'esecuzione sequenziale arriva a richiedere mesi se non anni di computazione.

E' naturale quindi pensare a Word Count come un problema di programmazione distribuita, in cui la distribuzione dell'input tra i nodi di un cluser può significa-

tivamente diminuire il tempo necessario, passando da anni al tempo per un caffè
:coffee:

2. Soluzione proposta

Vieni quindi proposta una soluzione al problema *Word Count* in ambiente distribuito, dove più nodi hanno accesso agli stessi dati. L'input è un insieme di file, e l'output è la lista delle parole contenute nei file associate al numero di loro occorrenze.

La soluzione proposta vuole massimizzare i vantaggi derivanti dalla distribuzione del calcolo e dalla *data locality*, massimizzando il tempo di computazione dei nodi e minimizzando quello di comunicazione. Il protocollo di comunicazione utilizzato per la comunicazione e distribuzione del calcolo è Message Passing Interface (MPI) ¹, il quale rappresenta lo standard *de facto* nel campo della comunicazione tra nodi in un cluster, di cui è stata utilizzata l'implementazione OpenMPI ²



Figure 2: OpenMPI logo

La soluzione è stata scritta in linguaggio C, e ne viene in seguito presentata e analizzata l'architettura che la compone.

2.1. Architettura

La soluzione segue il paradigma **Map-Reduce**, in cui la computazione è composta da una fase di *map*, in cui ogni nodo processa una parte dell'input, e da una fase di *reduce*, in cui i risultati dei singoli nodi vengono messi insieme.

L'input è fornito a tutti i nodi ed è composto da una path in cui sono presenti i file da esaminare. **Si assume che tutti i nodi abbiano accesso ai file alla path di input**

Il processo 0 provvederà a dividere il *workload* tra i nodi del cluster ed inviare loro i dati necessari (incluso se stesso), quindi ogni nodo processa la sua porzione di input. I dati vengono raccolti dal processo 0, che si occuperà di fornirli in output.

¹<https://www.mpi-forum.org/docs/mpi-3.1/mpi31-report.pdf>

²<https://www.open-mpi.org/>

La soluzione presentata è pertanto divisibile in 3 fasi:

1. Partizionamento e distribuzione dell'input ai nodi;
2. *Word Count* della porzione assegnata (**Map**);
3. Collezionamento dei risultati (**Reduce**).

2.1.1. Partizionamento e distribuzione dell'input

In questa fase è solo il processo 0 ad eseguirla, in quanto è propedeutica alla fase di calcolo distribuita.

Lo scopo di questa fase è quello di definire un partizionamento equilibrato del calcolo tra i vari nodi, al fine di minimizzare i tempi di *idle* derivanti da una distribuzione sbilanciata.

L'approccio scelto è **byte level**: ad ogni nodo è assegnata una porzione del file di una taglia fissata di byte, uguale per ogni nodo. Le motivazioni dietro questa scelta sono chiare: - Distribuzione equa: i nodi processeranno un quantitativo di byte uguale tra loro, uniformando il numero di accessi in lettura ai file necessari; - Disaccoppiamento tra i nodi: ogni nodo riceve un numero di byte da processare ed il file da cui partire, e computa finchè ha da leggere senza interessarsi dello stato degli altri nodi. Ogni nodo riceve quindi numero di byte da leggere ed il file di partenza: questo perchè i byte potrebbero attraversare più file. E' imperativo quindi che questi siano ordinati in egual modo per ogni nodo

Andiamo ad analizzare alcuni passi critici nell'esecuzione:

Per ogni nodo, il nodo 0 calcola le informazioni necessarie: nella struttura *info* è presenta il numero di byte da leggere, il file di partenza ed il byte di partenza. Il nodo assegna numero di byte da leggere uguale o di 1 unità maggiore a quello degli altri nodi (l'unità in più è necessaria per la gestione delle dimensioni che non dividono il numero di nodi) ed il file a lui assegnato.

Per calcolare il file a cui il nodo *i* arriva al termine della sua computazione (il quale sarà l'input per il nodo *i+1*) somma al byte di partenza del nodo *i*-esimo il numero di byte da lui letti e gli sottrae iterativamente le dimensioni del file attuale e successivi fino a che non trova il file di arrivo (in cui **start_byte** è più piccolo della sua dimensione). Se però **start_byte** è esattamente uguale alla dimensione del file corrente, allora significa che il nodo *i* legge esattamente fino alla sua fine e pertanto il nodo *i+1* inizia a leggere dal successivo.

```
1 for (int i = 0; i < num_proc; i++)
2     {
3         partitioning[i] = overflow && overflow_value-- > 0 ?
          total_dim /
4         num_proc + 1 : total_dim / num_proc;
5         infos[i].num_file = current_file;
6         infos[i].num_byte = start_byte;
7
8         if (i != (num_proc - 1))
9             {
```

```

10     start_byte += partitioning[i];
11     for (; current_file < num_files; current_file++)
12     {
13         if (start_byte <= dimensions[current_file])
14         {
15             if (start_byte ==
16                 dimensions[current_file])
17             {
18                 start_byte = 0;
19                 current_file++;
20             }
21             break;
22         }
23         start_byte -= dimensions[current_file];

```

A questo punto è necessario calcolare il byte di partenza del nodo $i+1$: a questo punto **start_byte** si trova all'ultimo byte letto dal nodo i , ma questo potrebbe trovarsi nel mezzo di una parola. L'algoritmo sceglie quindi di portare **start_byte** alla posizione della prima parola che il nodo $i+1$ dovrà leggere, che non è altro che la parola successiva a quella tagliata.

Per farlo saltiamo al byte **start_byte-1** (solo se non partiamo dall'inizio) e consumiamo ogni carattere che non sia un delimitatore di parola. Dato che potremmo terminare il file nel caso in cui la parola tagliata sia l'ultima, tale condizione viene ulteriormente verificata.

```

1      //Need to check if the current byte cuts a word
2      sprintf(path, "%s/%s", basepath,
3              files[current_file]->d_name);
4      f = fopen(path, "r");
5
6      if(start_byte != 0){
7          fseek(f, start_byte - 1, SEEK_SET);
8          while (!is_delimiter(c = fgetc(f)))
9          {
10             start_byte++;
11             partitioning[i]++;
12         }
13
14         //The cutted word could be the last one
15         if (feof(f))
16         {
17             start_byte = 0;
18             current_file++;
19         }

```

```

20     }
21
22     infos[i].counting = partitioning[i];
23 }

```

Il lettore potrebbe domandarsi il perchè non sia stato utilizzato un approccio *word level* o *line level* che eliminerebbe molta della complessità presenti, ed il motivo è presto spiegato ed è figlio del tipo di input a disposizione: entrambi gli approcci, seppur equilibrando al meglio il grado di distribuzione dell'input tra i nodi (ma non gli accessi al file) necessitano di una fase di *pre-processing* iniziale in cui contare il numero di parole o linee presenti nei file. Tale fase può essere estremamente lunga per file di dimensioni considerevoli, e pertanto poco performante.

2.1.2. Word count

Per il conteggio delle parole si utilizza una struttura **cell** composta da chiave-valore così definita:

```

1 typedef struct c_cell{
2     char key[30];
3     long int value;
4 } cell;

```

la dimensione *30* deriva dalla lunghezza massima delle parole nell'alfabeto italiano, che è 29 per la parola “*esofagodermatodigiunoplastica*”.³

Per la gestione delle parole in memoria è stata utilizzata una **hash map** implementata *ad hoc*, basa sull'algoritmo di hashing FNV su un massimo di 100000 liste.

Questo valore deriva da una precisa considerazione: questa implementazione vuole risolvere il problema del *Word Count* su documenti *reali*, cioè scritti in un linguaggio umano e che rispettino nelle forme una sintassi ben definita. Ne consegue che il numero di parole considerabili è definibile in funzione del numero di parole presenti (ed utilizzate) in tale linguaggio. In questo caso si è scelto di prendere come riferimento la lingua italiana, la quale conta circa 50000 parole di utilizzo comune⁴, valore raddoppiato per ridurre le collisioni in caso di presenza di termini ricercati o dialettali. Il numero di liste è comunque modificabile tramite il metodo *initialize_map*, che in questa soluzione è definito come *num_bytes / MAX_WORD_LENGTH*

```

1 static unsigned hashing(char *key)
2 {
3     for (int i = 0; i < strlen(key); i++)
4         key[i] = tolower(key[i]);

```

³<https://www.focus.it/cultura/curiosita/qual-e-la-parola-piu-lunga-della-lingua-italiana>

⁴http://www.treccani.it/magazine/lingua_italiana/domande_e_risposte/lessico/lessico_267.html

```

5
6     unsigned char *p = key;
7     unsigned len = strlen(key);
8     unsigned h = 2166136261;
9     int i;
10
11     for (i = 0; i < len; i++)
12     {
13         h = (h * 16777619) ^ p[i];
14     }
15
16     return h % num_lists < 0 ? h % num_lists * -1 : h %
        num_lists;
17 }

```

Fowler–Noll–Vo hash: la presenza del valore *-1* nel return è dovuto ai caratteri wchar, che nel prodotto forniscono valori negativi

Per la lettura delle parole nel file si fa uso di un modulo tokenizzatore, il quale legge da un file e restituisce la prossima parola delimitata da un carattere di delimitazione.

Un **delimitatore** è un carattere speciale che identifica la fine di una parola: uno spazio, un newline, una tabulazione, della punteggiatura etc. sono considerati delimitatori. L'elenco dei delimitatori utilizzati in questa soluzione è presente nel file "delimiters.txt".

E' possibile fornire anche un proprio file di delimitatori.

La fase di conteggio è quindi effettuata da ogni nodo, avendo come unica metrica il byte di partenza, il file di partenza ed il numero di byte da leggere.

Questa risulta molto semplice, in quanto non bisogna far altro che leggere una parola dal file ed inserire nella hash map una coppia <parola,1> che, se presente, aggiornerà il valore precedente nella hash map oppure, se non presente, creerà un nuovo nodo per quella parola. Quindi viene decrementato il numero di byte da leggere e viene cambiato il file se necessario.

```

1     initialize_map(bytes_to_read / MAX_WORD_LENHT);
2     new_hash_map();
3     while (bytes_to_read > 0 && current_file < num_files)
4     {
5         if (file == NULL)
6         {
7             sprintf(path, "%s/%s", basepath,
                files[current_file]-> d_name);
8             file = fopen(path, "r");
9             if (start_byte > 0) //Skip to start position
10            {
11                fseek(file, start_byte, SEEK_SET);

```



```

12         start_byte = 0;
13     }
14 }
15
16     readed_bytes = next_word(file, &word);
17     bytes_to_read -= readed_bytes;
18     if (word == NULL)
19     {
20         fclose(file);
21         current_file++;
22         file = NULL;
23         continue;
24     }
25
26     if (bytes_to_read >= 0)
27     {
28         //Actually add the word to the hash map
29         strcpy(word_cell.key, word);
30         add_or_update(&word_cell);
31     }
32
33     free(word);
34 }

```

2.1.3. Reduce

La fase di reduce prevede il collezionamento dei risultati dei vari nodi verso il nodo 0, il quale si occuperà dell'output finale.

La struttura adottata è ad **albero rovesciato direzionato**: - I nodi rappresentano le foglie dell'albero

- Ogni arco rappresenta:
- Invio dei risultati se diretto ad un nodo con etichetta diversa dal nodo di partenza
- Ricezione dei risultati dai nodi origine degli archi entranti
- Il nodo radice è il nodo obiettivo della reduce

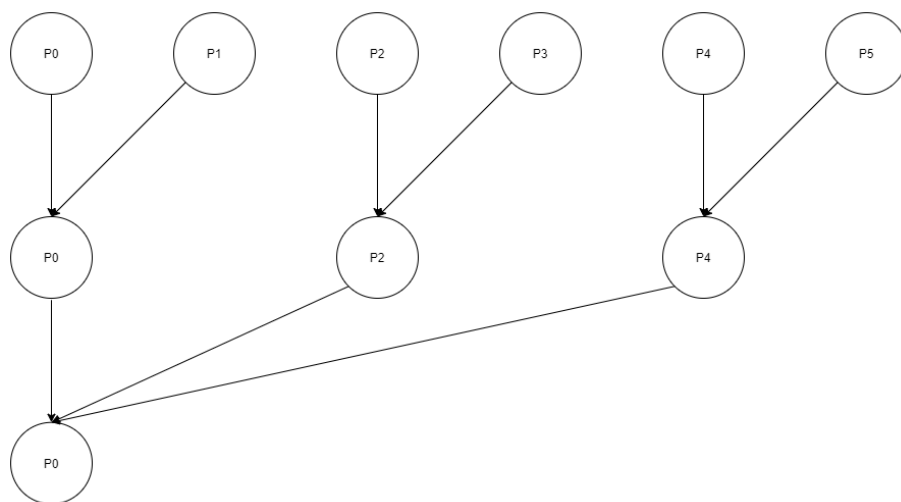


Figure 3: Reduce

La scelta di tale struttura è dovuta alla volontà di massimizzare l'utilizzo dei nodi nella fase di reduce e distribuire di conseguenza il calcolo in modo intelligente tra i nodi del cluster: infatti tale scelta permette di eseguire, per n nodi, al più $\lceil \log_2 n \rceil$ *receive* e ridurre considerevolmente il carico di lavoro che, altrimenti, sarebbe gravato su un singolo nodo.

La comunicazione in questa fase è essenzialmente *asincrona*: ogni nodo definisce dei buffer di ricezione per ogni nodo da cui dovrà ricevere dati che verranno riempiti asincronamente. Ciò ha il vantaggio di massimizzare la velocità di ricezione al costo di un incremento nell'uso di memoria RAM, evitando inoltre stalli nella fase di ricezione a causa di nodi più lenti di altri.

L'algoritmo può essere scomposto in 3 fasi:

1. Calcolo di *receivers* e *sender*;
2. Ricezione dei dati
3. Invio dei dati

2.1.3.1. Calcolo *receivers* e *senders* Ogni nodo ha necessità di calcolare a priori quali saranno i nodi da cui riceverà i dati temporanei ed a quale nodo dovrà inviare i dati da lui computati. Questo è necessario per sfruttare al meglio i sistemi di comunicazione asincrona.

Per fare ciò **ogni nodo attraversa l'albero e calcola il suo comportamento per ogni livello:**

Di default, in ogni livello i nodi in posizione dispari inviano dati al nodo pari alla loro sinistra, mentre i nodi pari ricevono dal nodo a destra. Per ogni livello ogni nodo calcola quindi la propria posizione, da chi dovrà ricevere, per il quale

alloca un buffer apposito e lancia una routine di ricezione asincrona *MPI_Irecv*, ed aggiorna il nodo a cui dovrà inviare i risultati.

Eccezione va fatta per i nodi pari ultimi: se infatti questi sono potenze del 2, allora comunicano direttamente col nodo 0 in quanto non riceveranno mai dati, altrimenti continuano la computazione al livello successivo.

```

1 send_proc = rank - 1; //Default behaviour
2   while (cur_rank % 2 == 0 && level < ceil(log2(num_proc)))
3   {
4       //I'm the last one
5       if (ceil(num_proc / pow(2, level)) == (cur_rank + 1)){
6           //If I am a power of two and last one, i communicate
7           with rank 0
8           if (is_power_of_two(cur_rank))
9           {
10              send_proc = 0;
11              break;
12          }
13      }
14      else
15      {
16          if (num_requests == dim_buffer)
17          {
18              dim_buffer *= 2;
19              buffers = realloc(buffers, sizeof(cell*) *
20                              dim_buffer);
21          }
22          buffers[num_requests] =
23              calloc(MAX_WORD_NUM, sizeof(cell));
24          MPI_Irecv(buffers[num_requests], MAX_WORD_NUM *
25                    sizeof(cell), CELL,
26                    (int)(rank + pow(2, level)),
27                    0, MPI_COMM_WORLD,
28                    &requests[num_requests]);
29          num_requests++;
30      }
31      cur_rank = cur_rank / 2;
32      level++;
33      send_proc = (int)(rank - pow(2, level));
34  }

```

2.1.3.2. Ricezione dei dati Una volta calcolate tutte le informazioni necessarie, ogni nodo aspetta che tutte le routine di ricezione terminino.

Qui diventa evidente il vantaggio dell'asincronia: ogni nodo infatti aspetta semplicemente di ricevere il segnale che una ricezione è completata tramite la

funzione *MPI_Waitany* e la computa, disinteressandosi dell'etichetta del nodo sorgente e del livello al quale è stata definita tale routine. Ciò permette di evitare stalli e rallentamenti nella ricezione dovuti a nodi meno performanti, i quali avranno più tempo per computare mentre altri nodi più veloci mantengono occupato il nodo destinazione.

In sintesi, tale sistema può considerarsi puramente meritocratico.

```

1 int index,received_dim;
2     MPI_Status status;
3     for (int i = 0; i < num_requests; i++)
4     {
5         MPI_Waitany(num_requests,requests,&index,&status);
6         printf("RANK %d: received data from rank
7               %d\n",rank,status.MPI_SOURCE);
8         MPI_Get_count(&status,CELL,&received_dim);
9         add_cell_array(buffers[index],received_dim);
10        free(buffers[index]);
11    }
12    free(buffers);

```

2.1.3.3. Invio dei risultati Questa fase consiste semplicemente nell'invio dei propri risultati parziali al nodo di destinazione precedentemente calcolato.

Unica nota è la funzione *compact_map_ordered()*:

```

1     long int dim;
2     cell *results = compact_map_ordered(&dim);

```

Tale funzione restituisce un vettore di strutture *cell* ordinate rispetto alla chiave (cioè alla parola), ed è necessario per minimizzare le primitive di comunicazione necessarie. La primitiva di invio *MPI_Send*, infatti, ha un funzionamento simile alla funzione *write* del linguaggio C, pertanto richiede un indirizzo base ed il numero di locazioni di memoria consecutive da inviare. Ciò rende necessario il compattamento della struttura *hash map* utilizzata, la quale per sua natura non è contigua in memoria.

3. Benchmark

La soluzione proposta è stata testata apporfonditamente al fine di verificarne correttezza e scalabilità in un ambiente distribuito, facendo attenzione al definire delle modalità di testing il più oneste possibili. Ciò significa utilizzare input e casi d'uso che possano rispecchiare quanto più possibile un reale input ad un tale sistema in ambienti di produzione.

Le principali metriche utilizzate per la valutazione delle prestaizoni sono lo *Speedup* e l' *Efficienza*:

Lo *Speedup* indica l'incremento prestazionale sullo stesso input tra un'esecuzione sequenziale ed una distribuita, calcolato con la seguente formula:

$$Speedup = \frac{Time_{sequential}}{Time_{parallel}}$$

Per n processori, lo speedup massimo è ovviamente n .

Tale parametro ci da un'utilile parametro sulla bontà della nostra soluzione.

L'*Efficienza* è una normalizzazione dello *Speedup* e definisce la bontà della soluzione indicando quanto i tempi di esecuzione distribuita e sequenziale si avvicinino:

$$Efficienza = \frac{Time_{sequential}}{n * Time_{parallel}}$$

Per l'efficienza, il limite superiore è 1 e rappresenta il risultato migliore auspicabile.

3.1. Architettura di testing

Il testing è stato effettuato sulla piattaforma Amazon AWS, in cui è stato configurato un cluster di 16 istanze EC2 modello m4.large.

Le istanze m4 rappresentano una scelta bilanciata in memoria, calcolo e rete ed il loro equilibrio è ideale per realizzare un ambiente di testing efficiente e costante. > La scelta del modello, rispetto alle più economiche T2, deriva dal *CPU Credit* presente su queste: ogni credito permette di ottenere prestazioni maggiori per 1 minuto di tempo, ed è pensato per applicazioni web che presentano picchi di carico nel corso della giornata.

> Nel nostro ambiente di testing tale variabile ha portato a risultati altalenanti, pertanto si è preferito evitare completamente tale modello.

Ogni macchina m4.large presenta la seguente configurazione hardware:

- Processori Intel Xeon® E5-2686 v4 (Broadwell) da 2,3 GHz (2 vCPU);
- 8GB RAM;
- 450 Mb\s bandwidth
- 25 GB storage EBS

La configurazione software è invece composta da Ubuntu® Bionic 18.04, con installati solo una piccola serie di software aggiuntivi quali: - Htop - VIM > La configurazione è possibile riprodurla configurando una macchina EC2 con AMI *ami-07ebfd5b3428b6f4d*

Durante la fase di testing le macchine erano scariche di qualsiasi processo aggiuntivo se non quelli previsti di base dal sistema operativo, ad eccezione di una sessione SSH aperta sul nodo 0 ad inizio fase di test, ma assolutamente irrilevante nel quadro delle prestazioni.

3.2. Esperimenti

Gli esperimenti effettuati sono stati definiti al fine di verificare la capacità della soluzione proposta di *scalare* in un ambiente distribuito reale.

Per *scalabilità* si intende la capacità di un applicativo software di migliorare le proprie performance al crescere delle risorse disponibili sulla macchina corrente (*scalabilità verticale*) o all'aumentare dei nodi che partecipano alla computazione (*scalabilità orizzontale*). Nei nostri esperimenti prenderemo in considerazione solo la scalabilità orizzontale, la quale meglio definisce la bontà della soluzione proposta in un ambiente distribuito.

Tale metrica viene ulteriormente divisa in 2 metriche differenti:

- *Strong Scalability*
- *Weak Scalability*

Presenteremo tali metriche singolarmente, associate ai risultati ottenuti.

3.2.2. Strong Scalability

Nella *Strong Scalability* la taglia dell'input resta costante mentre il numero di processori aumenta. Tale metrica permette di definire quanto bene l'applicativo è in grado di velocizzare il calcolo di una certa istanza del problema all'aumentare dei nodi: maggiore è l'incremento, maggiore è la qualità della soluzione.

La *Strong Scalability Efficiency* è calcolata tramite la seguente formula:

Sia n il numero di processori, t_1 il tempo di esecuzione per 1 processore e t_n il tempo per n processori, allora la Strong Scalability Efficiency è :

$$E_{strong} = \left(\frac{t_1}{n * t_n} \right)$$

Per questi esperimenti è stato utilizzato come input un insieme di file aventi 1 parola per linea per una taglia totale di 6 GB.

Qui sono riportati i risultati ottenuti:

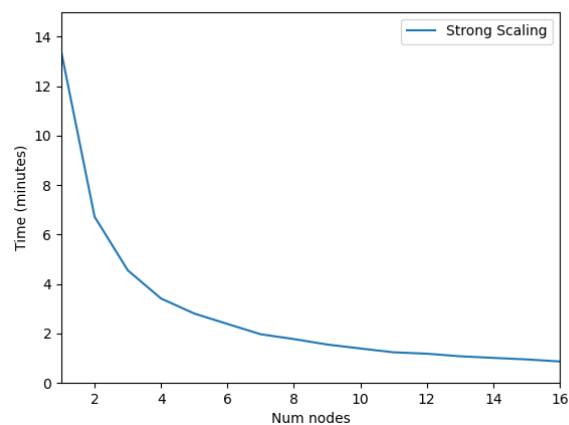


Figure 4: Tempo di esecuzione per scalabilità forte

Come è possibile osservare dalla figura 3, l'efficienza si avvicina al valore 1 per tutti gli esperimenti effettuati, il che indica un'ottimo scaling da parte della soluzione proposta.

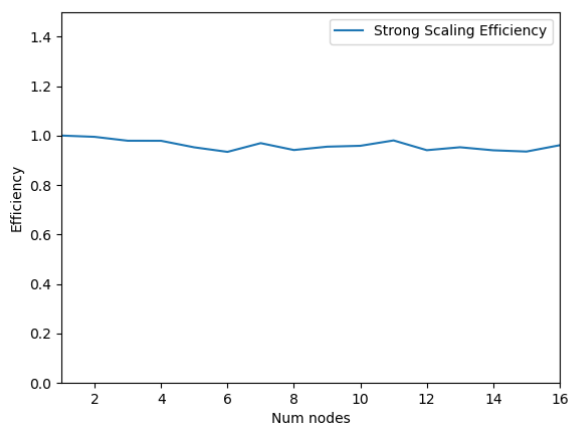


Figure 5: Efficienza per scalabilità forte

3.2.1. Weak Scalability

Nella **Weak Scalability** si verificano le prestazioni di un applicativo software quando la taglia dell'input cresce proporzionalmente al numero di processori. Tale metrica misura in particolar modo l'impatto dell'overhead derivante dalla comunicazione nell'ambiente distribuito sulle performance dell'applicazione.

La **Weak Scalability efficiency** è calcolata tramite la seguente formula:

Sia n il numero di processori, t_1 il tempo di esecuzione per 1 processore e t_n il tempo per n processori, allora la Weak Scalability Efficiency è :

$$E_{weak} = \frac{t_1}{t_n}$$

Per questi esperimenti è stato utilizzato un file contenente 1 parola per riga dalla taglia di *1GB*, replicato per ogni nodo partecipante alla computazione. Per esempio, per 3 nodi il file era replicato 3 volte per una taglia totale di input di *3GB*

Qui sono riportati i risultati ottenuti:

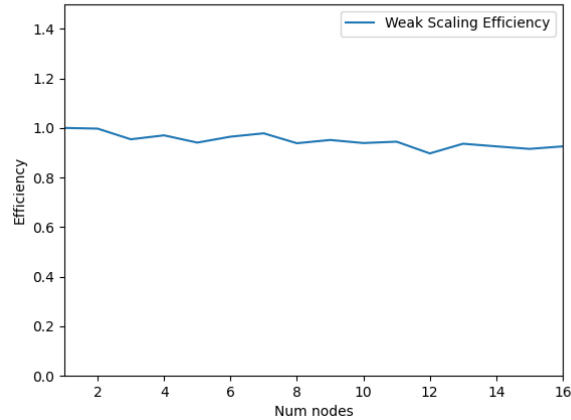


Figure 6: Efficienza per scalabilità debole

Come è possibile osservare dalla figura 4, l'andamento della soluzione proposta peggiora leggermente all'aumentare dei nodi: ciò è naturale ed è dovuto all'overhead sempre maggiore generato dalla comunicazione dei risultati intermedi da parte dei nodi. Ciononostante, l'efficienza resta sempre molto vicino al valore 1 ottimale, dimostrazione di una strutturazione efficiente.

4. Conclusioni

Abbiamo presentato il problema del Word Count, che consiste nel determinare il numero di occorrenze di ogni parola presente in un insieme di file. Abbiamo quindi presentato una soluzione distribuita del problema, implementata utilizzando lo standard di comunicazione MPI e ne abbiamo analizzato accuratamente le performance.

La soluzione proposta ha ottenuto ottimi risultati sia in scalabilità forte che debole, dimostrando stabilità e prestazioni eccellenti.

Possibile miglioramento potrebbe essere quello di definire un miglior fattore di determinazione del numero di parole aspettate nei testi da analizzare, utilizzando tecniche probabilistiche, che potrebbe decrementare significativamente il consumo di memoria per ogni operazione di *receive* di ogni nodo.