

# MINHASH PROBLEM

# INTRODUZIONE AL MINHASH

- MinHash è una tecnica algoritmica per stimare velocemente il grado di somiglianza tra due insiemi.
- Solitamente è applicato su larga scala in grandi problemi di raggruppamenti di documenti per somiglianza dei loro insiemi di parole.
- Inizialmente usato dal motore di ricerca AltaVista per individuare duplicati ed eliminarli dai risultati ricerca.
- Implementato da Andrei Broder nel 1997.

# LA JACCARD SIMILARITY

Per il calcolo della somiglianza tra insiemi si usa comunemente il coefficiente di similarità di Jaccard, che è definito dal rapporto tra la dimensione dell'intersezione e la dimensione dell'unione degli insiemi campionari

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}.$$

# OBIETTIVO MINHASH

- L'obiettivo del MinHash è di stimare velocemente  $J(A,B)$  senza calcolare esplicitamente l'intersezione e l'unione.
- Tramite una funzione di hash e delle permutazioni applicate a insiemi di dati, di cui vengono ricavati dei valori di hash minimi, il minHash produce una stima (probabilistica) approssimativa dell'Indice di Jaccard.

# I TASK INDIVIDUATI

- “List dir” recupero da un path dato dei path dei files da confrontare.
- “tokenizer” recupero del contenuto di un file e pulizia del testo da caratteri spurii.
- “Shingle extract” creazione degli shingles dato un testo
- “Get Signatures” creazione delle signature tramite permutazioni di hash dati gli shingles.
- “Get Similarities” tramite un’opportuna comparazione vengono stampate e individuate le somiglianze.

# ALGORITMI PARALLELI UTILIZZATI

- OMP e parallelismo sui dati
- Pthreads e parallelismo sui dati
- Approccio con l'utilizzo di tutte e due le librerie assieme e parallelismo basato principalmente sui tasks.

# OMP

- Propensione della libreria a parallelizzare algoritmi seriali già diviso in piccoli tasks.
- Parallelizzazione sui dati
- Utilizzo dell'istruzione parallel e ricerca delle critical section
- Utilizzo delle istruzione reduce e atomic per risolvere le problematiche dovute alla parallelizzazione.

```

graph TD
    main([main  
prende in input il path  
di una directory]) --> list_dir_parallel[omp parallel for  
list_dir()  
list_dir()  
...  
list_dir()]
    list_dir_parallel --> minHash_parallel[omp parallel for  
minHash()  
...  
minHash()]
    minHash_parallel --> string_cleaned_parallel[omp parallel for  
string_cleaned()  
...  
string_cleaned()]
    string_cleaned_parallel --> shingle_extract_parallel[omp parallel for  
shingle_extract()  
...  
shingle_extract()]
    shingle_extract_parallel --> get_signatures_parallel[omp parallel for  
get_signatures()  
...  
get_signatures()]
    get_signatures_parallel --> find_similarities[find_similarities()]
    find_similarities --> get_sketches_parallel[omp parallel for  
get_sketches()  
...  
get_sketches()]
    get_sketches_parallel --> merge_sort_parallel[omp parallel for  
merge_sort()  
...  
merge_sort()]
    merge_sort_parallel --> create_triplets_parallel[omp parallel for  
create_triplets()  
...  
create_triplets()]
    create_triplets_parallel --> do_clustering_parallel[omp parallel for  
do_clustering()  
...  
do_clustering()]
    do_clustering_parallel --> check_similarity[check_similarity()]
    
```



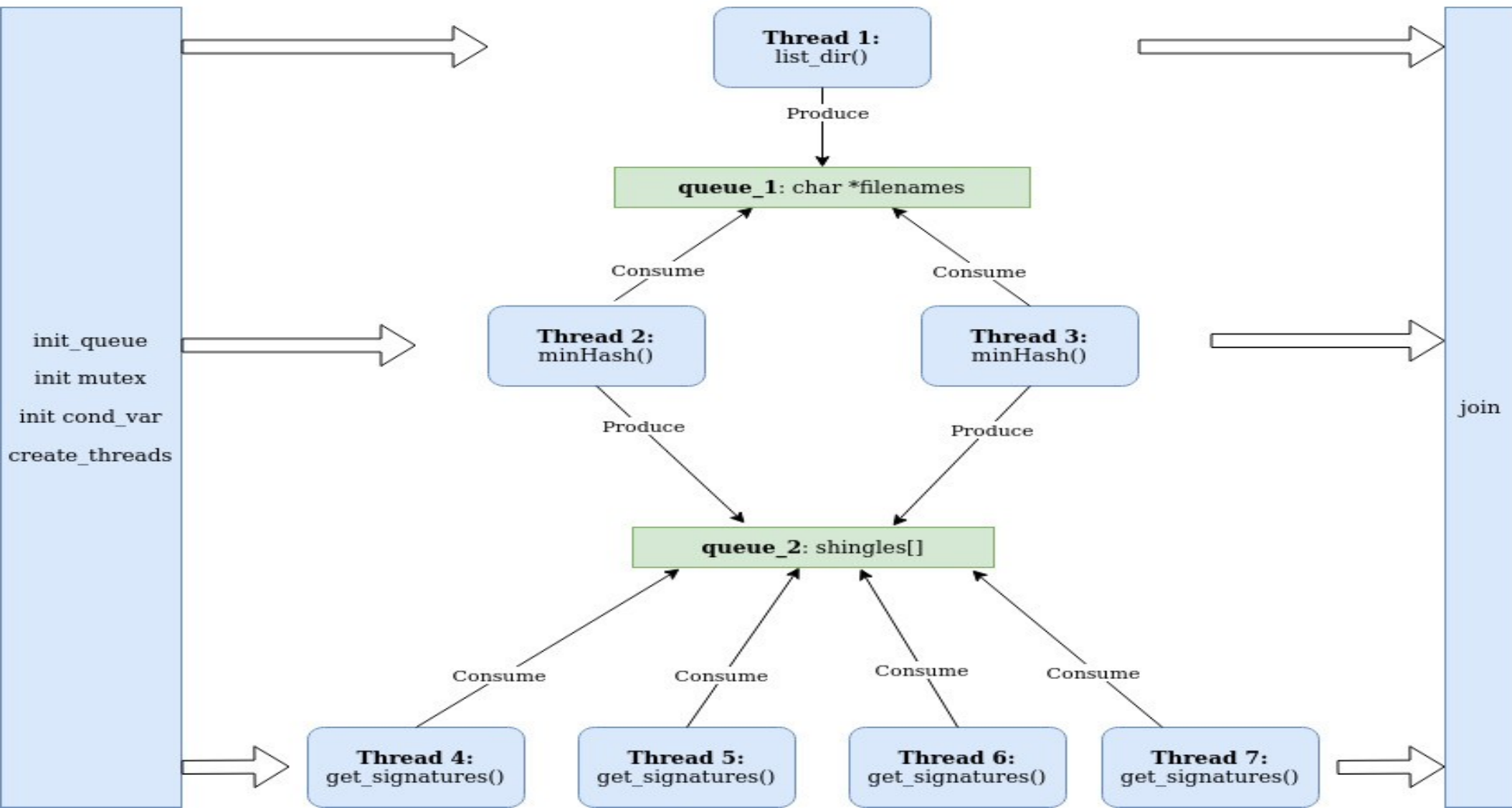
# PTHREAD

- Parallelizzazione sui dati
- Approccio più a basso livello rispetto a OMP con uso di strutture dati apposite per comunicare con i thread creati.
- Individuazione delle barrier (`pthread_join`)
- Uso di mutex e individuazione delle sezioni critiche

# Produttore Consumatore

- In questa versione è stata utilizzata principalmente la libreria pthread, solo nel task più costoso si è deciso di usare omp.
- Visione dei task come Produttori e Consumatori di strutture dati definite
- Parallelizzazione sui Task
- Struttura dati della coda per comunicare tra task

## Main

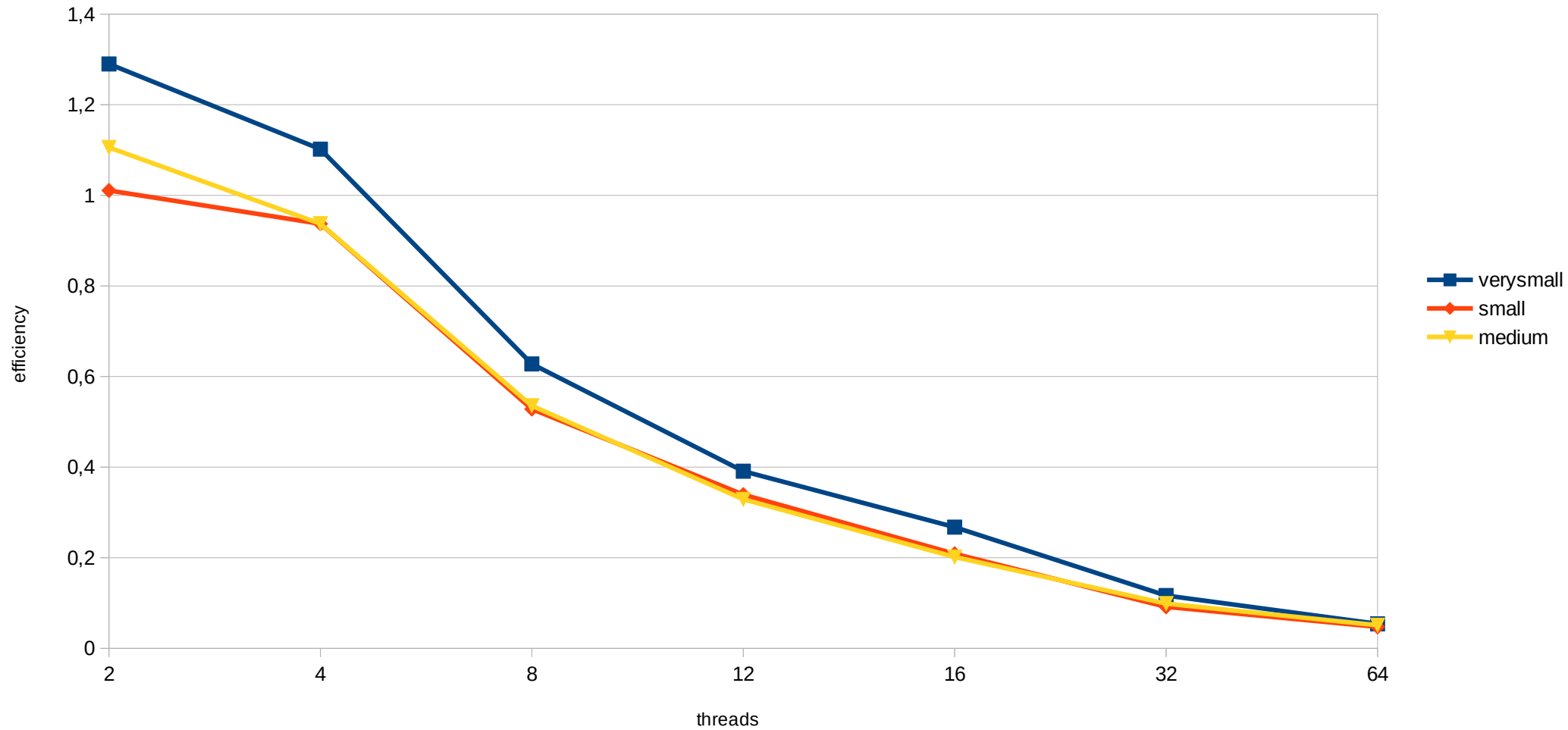


# TEST

- Raccolta dati basata sulle misurazioni dei wall time
- Controllo degli errori sulle signature create nella versione seriale
- Speed-up
- Efficienza e scalabilità

# OpenMP efficiency on threads

no internal threads



## OpenMP efficiency on threads

4 nested threads (Esempio 8 = 2 main threads \* 4 nested threads)

