

Programový model nad distribuovanou pamětí : MPI

(procesy, komunikátory, 2-bodové ^{skupinové} a komunikační operace, blokující a neblokující operace a jejich komunikační módy)

Hybridní MPI + OpenMP model

MPI - Message passing interface → rozhraní zpráv mezi procesy paralelního programu

- pracuje nad distribuovanou pamětí (NUMA)
- každý proces je vždy součástí alespoň jedné skupiny (číslování od 0 do # - 1)
 - v každé skupině obdrží jiné číslo

Komunikátory - součástí (parametr) každé komunikační operace

- určuje v rámci které skupiny probíhá komunikace - MPI_COMM_WORLD pro všechny procesy
- index (uvnitř skupiny) a index (mezi skupinami)

MPI_Comm_rank - číslo procesu v dané skupině

MPI_Comm_size - počet procesů dané skupiny

Komunikační operace

- o 2 bodové - mezi dvěma procesy
- o kolektivní - všechny s daným komunikačním
- o blokující - čeká se na splnění určité podmínky
- o neblokující - konec komunikace - nutné obsloužit později

MPI_Send (*buff, count, MPI_Datatype, dest, tag, comm)
 MPI_Recv (*buff, count, MPI_Datatype, src, tag, comm, &status)

↑ cílový proces
 ↑ zdroj proces
 ↑ typ daných přenášených dat [MPI_INT, MPI_DOUBLE, ...]
 ↑ značka přenášených dat [MPI_ANY_TAG]
 ↑ od koho mám přijímat [MPI_ANY_SOURCE]

struktura obsahující source a tag příchozí zpráv [MPI_STATUS_IGNORE]

- nezákladní blokující operace (= buffered nebo synchronous mode)

MPI_Bsend - buffered → uložení do systémového bufferu (nečekaně na cíli)

MPI_Ssend - synchronous → končí až když cílový proces iniciální příjem

MPI_Rsend - ready mode → příjem ne musí být iniciován, jinak konec s chybou

- všechny operace mají i neblokující variantu (MPI_Isend , MPI_Irecv , ...)
 - musíme explicitně řešovat dokončení operace
 - buffer do té doby nelze modifikovat
- funkce mají dodatečný parametr MPI_Request
 - MPI_Test , MPI_Wait & ... any, ... all pro hromadění operací
- speciální MPI_Sendrecv
- MPI_Probe - kontroluje příchod zpráv

Kolektivní komunikační operace

- 1:N → one-to-all broadcast - MPI_Bcast (stejný multicast)
- one-to-all scatter - MPI_Scatter
- all-to-one gather - MPI_Gather - shromáždění dat ze všech
- poslední dvě se liší pouze ve směru, takže jsou stejné

N:N

- all-to-all broadcast
- = all-to-all gather - MPI_Allgather
- all-to-all scatter - MPI_Alltoall

Hybridní model MPI + OpenMP

- v rámci jednoho výpočetního uzlu / procesoru běží jeden nebo více MPI procesů, každý se dělí na několik vláken pomocí OpenMP
- 1. proces na jednom výpočetním uzlu → celý uzel na celý proces
- 2. proces na každém procesoru (tj. socket) → lepší přístup ke sdílené paměti
- musí se inicializovat pomocí MPI_Init_thread na požadovanou míru spolupráce
 - MPI_THREAD_SINGLE - dělíme se na vlákna
 - $\text{MPI_THREAD_FUNNELED}$ - pouze hlavní vlákno může volat MPI funkce
 - $\text{MPI_THREAD_SERIALIZED}$ - MPI funkce musí v jednom chvilku volat jen jedno vlákno
 - $\text{MPI_THREAD_MULTIPLE}$ - všeobecný model