

Interaction overhead

① 줄이기 위해 고려할 점

- data locality ↑ ... Communication for data ↓
- graph partitioning ... crossing edge ↓ ... 연결 끊어 해결하기
- communication frequency ↓ ... 한번에 chunk/batch communicate 해서 frequency 줄이기
- contention/hot-spot ↓ ... bottleneck 방지하기 위해 de-centralized dynamic mapping 사용

② 방법

· overlapping communication-computation

- single-thread : interrupt 해서 CPU idle 줄이기
prefetch 사용해서 latency 줄이기
- multi-thread : communication을 other computation과 overlapping

· replicate data

- redundancy를 조금 허용해서 communication 줄이기
- false-sharing도 줄일 수 있음 (local copy 유지)

· overlapping communications

- stepwise : small tasks의 chain의 pipeline
- naive . 그냥 연결하기
- pipelined naive : 전체 pipeline

Minimizing Interaction Overhead

Rules of Thumb	experience-based method
Maximize data locality	<ul style="list-style-type: none"> don't fetch data you already have (no redundant data) reconstruct computation to reuse data promptly <p>Communication is expensive! locality ↑ communication ↓ ↳ fix graphed task flow (locality maximization) & reuse</p>
Minimize communication volume	<ul style="list-style-type: none"> partition interaction graph to minimize edge crossings ★ minimize communication data size & frequency for control flow
Minimize communication frequency	<p>chunk size ↑ size ↑ reuse ↑</p> <ul style="list-style-type: none"> aggregate message where possible use batch (chunk) saturate memory bandwidth <p>→ all Best</p>
Minimize contention/hot-spot	<ul style="list-style-type: none"> no contention on use of same network link / memory block <ul style="list-style-type: none"> can occur bottleneck no specific process heavily loaded (hot-spot) <ul style="list-style-type: none"> can be load-imbalancing and idling decentralized techniques <ul style="list-style-type: none"> centralized can make bottleneck on master thread

Master + slave
distributed

Minimizing Interaction Overhead (Cont.)

Techniques

Overlap Communication - Computation

Single-thread

use non-blocking communication primitives

- 즉, I/O processing 동안 CPU computation 할 수 있도록 interrupt 같은 policy 중요!
- overlap communication with your own computation ... 2개까 context switch ok
- one-sided communication (P2W) : prefetch remote data to hide latency

Thread switch?

Multi-thread

• multithread code

: overlap communication with another thread's computation

replicate data/computation to reduce communication

: communication을 줄이기 위해서 각각이 모든 process가 데이터를 갖고있도록

- @ dense matrix multiplication의 vector b copy 하든??

peer-only는 괜찮지만 write data-by
inconsistency 발생

Use collective interaction instead peer-to-peer

: data transfer나 contention overhead를 줄일 수 있다

Issue Multiple communications

: reduces exposed latency. 즉, send를 한꺼번에 하거나 pipelining 할 수 있는가?

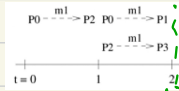
↓
필요 데이터 미리 받아
재발-하루!

Overlapping Communications

- P_0 broadcast data to others (P_1, P_2, P_3)
4 message

Step-wise algorithm

- Commonly used algorithm



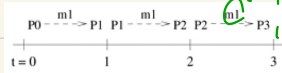
매각한을

매각한을

- takes 8 time to broadcasting

naive algorithm

- 2배로!



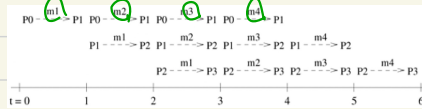
- takes 12 time for broadcasting

Pipelined naive algorithm

무엇까지든 끝내 된다면,
비교가동가능!

pipelined fashion

- 2배로!는 맞는데, pipelining 하시오!



- takes 6 time for broadcasting