# What Is Parallel Computing

Sequential task를 parallel task로 decompostion 하고
computer source를 각 tasks에 Mapping 시켜서
simultaneous operation을 통해 performance를 향상시키는 것

# Why parallel Computing

It can reduce execution time
It can efficiently use computer resource (without idle)

# Open MP

: open specification for Multi Processing
- execution model : fork-join
- memory model : shared memory

## false sharing → save cache line 공유
write on 것 -- sync 비용증가해서 performance 향상방해 것

Shared memory system 에서 빈번하게 performance reduction의 요인
when multiple thread occupy different core,
It shares L3 code and data is transferred in cache line unit.
when each thread access contiguous memory region,
Each element or chunk is smaller than cache line,
They share same cache line.
Then when they copy whole cache line which contains other elements.
If they do write operation on it,
Data Inconsistency / code coherence violated. so
It should synchronize cache line. validate(update occurs.

# Parallel Algorithm design

① two major steps ⇒ decomposition / Mapping

② critical path is longest weighted path
   It represents lowerbound/shortest time on that parallel algorithm

③

④ Task Interaction graph
   : data interaction이 중요함. weight edge = communication cost

⑤ recursive decomposition
   - when find minimum, we can decompose one serialized task
     into several chunks

① Data decomposition
   - Input / Input-output / Output / Intermediate data를
     decomposition 하는것.
   - data만 가르고 computation 그냥은 묵시적으로 대응임의 가정

① general design guide for minimizing interaction
   ① graph based static mapping, ~~output~~ cutting edge 최소로 가능총
   
   ② copy를 특히 interaction줄이기
   
   ③ communicate bandwidth 이상치초 volume↑ frequency↓
   
   ④ overlapping computation — communication