

Lecture 22: Cache-Oblivious Algorithms

Rafael Oliveira

University of Waterloo
Cheriton School of Computer Science

rafael.oliveira.teaching@gmail.com

November 30, 2020

Overview

- Administrivia
- Designing Algorithms in Real Life
- Cache-Aware and Cache-Oblivious Algorithms
- Conclusion
- Acknowledgements

Rate this course!

Please log in to

<https://evaluate.uwaterloo.ca/>

from *November 24th until December 7th* and provide us with your evaluation and feedback on the course!

- This would really help me figuring out what worked and what didn't for the course
- And whether I should put memes or gifs into my slides...
- Teaching this course is also a learning experience for me :)

Research Opportunities at UW!

Consider doing a URA, URF or USRA with a U Waterloo faculty!

See research openings at:

- Undergraduate Research Assistanship (URA):

[https://cs.uwaterloo.ca/computer-science/
current-undergraduate-students/research-opportunities/
undergraduate-research-assistantship-ura-program](https://cs.uwaterloo.ca/computer-science/current-undergraduate-students/research-opportunities/undergraduate-research-assistantship-ura-program)

- Undergraduate Research Fellowship (URF):

<https://grec.cs.uwaterloo.ca/>

- Undergraduate Research Internship (URI):

[https://cs.uwaterloo.ca/current-undergraduate-students/
research-opportunities/
undergraduate-research-internship-uri-program](https://cs.uwaterloo.ca/current-undergraduate-students/research-opportunities/undergraduate-research-internship-uri-program)

- For Canadians, please check out NSERC's USRA:

<https://cs.uwaterloo.ca/usra>

Memory Hierarchy

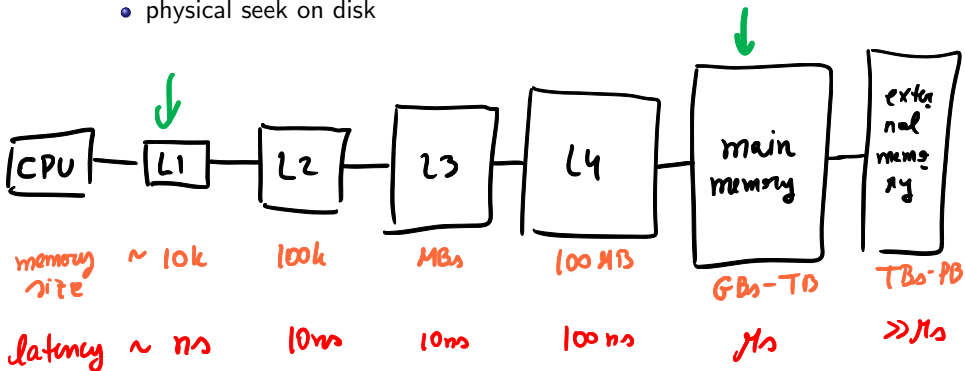
- Not all memory created equal

Memory Hierarchy

- Not all memory created equal
- In real life, need to account for our use of memory - it can be the bottleneck!

Memory Hierarchy

- Not all memory created equal
- In real life, need to account for our use of memory - it can be the bottleneck!
- Bigger memory usually have slower latency
 - distance travel to cache
 - physical seek on disk



Memory Hierarchy

- Not all memory created equal
- In real life, need to account for our use of memory - it can be the bottleneck!
- Bigger memory usually have slower latency
 - distance travel to cache
 - physical seek on disk
- Bandwidth (rate of transmission) is usually matched across memories

Memory Hierarchy

- Not all memory created equal
- In real life, need to account for our use of memory - it can be the bottleneck!
- Bigger memory usually have slower latency
 - distance travel to cache
 - physical seek on disk
- Bandwidth (rate of transmission) is usually matched across memories
- Memory architects *group* (or *block*) memory to decrease latency!
 - when fetching word of data, get entire *line/block* containing it
 - amortize latency over the whole line of memory

$$\frac{\text{latency}}{\text{block size}} + \frac{1}{\text{bandwidth}}$$

↑ ↑
set these to be equal
(adjusting block size usually)

Memory Hierarchy

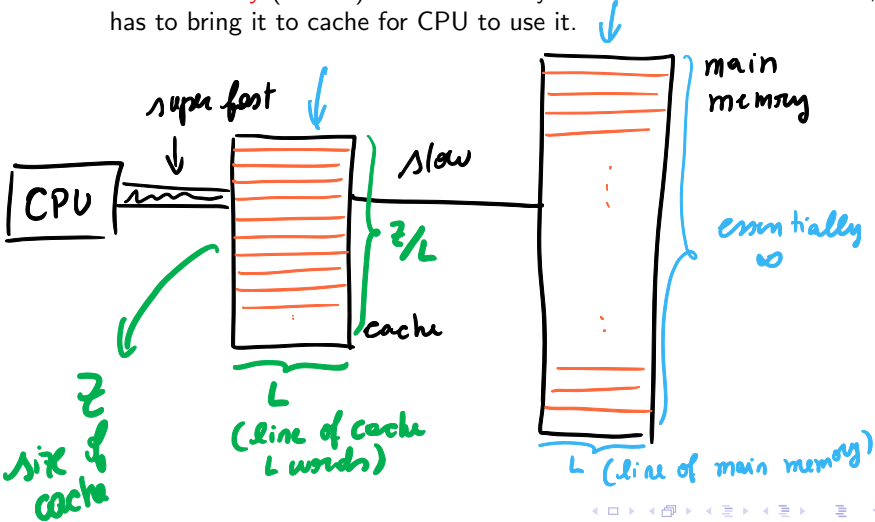
- Not all memory created equal
- In real life, need to account for our use of memory - it can be the bottleneck!
- Bigger memory usually have slower latency
 - distance travel to cache
 - physical seek on disk
- Bandwidth (rate of transmission) is usually matched across memories
- Memory architects *group* (or *block*) memory to decrease latency!
 - when fetching word of data, get entire *line/block* containing it
 - amortize latency over the whole line of memory

$$\frac{\text{latency}}{\text{block size}} + \frac{1}{\text{bandwidth}}$$

- Good algorithms must have:
 - *Spatial locality*: use all elements in same memory line
 - *Temporal locality*: re-use lines in cache before moving on

External Memory Model [Aggarwal and Vitter 1988]

- For simplicity, just 2 levels:
 - **Cache**: the fast memory (CPU has instant access to any word in it)
 - **Main memory** (or disk): slower memory. Once fetch data from there, has to bring it to cache for CPU to use it.



External Memory Model [Aggarwal and Vitter 1988]

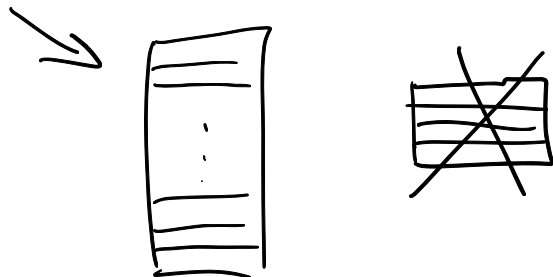
- For simplicity, just 2 levels:
 - *Cache*: the fast memory (CPU has instant access to any word in it)
 - *Main memory* (or disk): slower memory. Once fetch data from there, has to bring it to cache for CPU to use it.
- Algorithm has to read and write lines of memory!

External Memory Model [Aggarwal and Vitter 1988]

- For simplicity, just 2 levels:
 - *Cache*: the fast memory (CPU has instant access to any word in it)
 - *Main memory* (or disk): slower memory. Once fetch data from there, has to bring it to cache for CPU to use it.
- Algorithm has to read and write lines of memory!
- Complexity Measures:
 - *Work*: number of *computational* steps needed from start to end.
 - *Cache Misses*: number of cache misses (memory transfers from disk to cache) during entire computation

Cache-Oblivious Model [Frigo et al 1999]

- Ideal cache model:
 - word size known to everyone in advance
 - two-level memory
 - Memory partitioned in lines of L words each
 - Cache size is Z
 - *Tall Cache Assumption*: $Z = \Omega(L^2)$
 - Eviction algorithm: *farthest in the future*¹



¹Can be replaced by LRU efficiently.

Cache-Oblivious Model [Frigo et al 1999]

- Ideal cache model:
 - word size known to everyone in advance
 - two-level memory
 - Memory partitioned in lines of L words each
 - Cache size is Z
 - *Tall Cache Assumption*: $Z = \Omega(L^2)$
 - Eviction algorithm: *farthest in the future*¹
- *Cache-Oblivious Algorithm*: doesn't know Z, L

¹Can be replaced by LRU efficiently.

Cache-Oblivious Model [Frigo et al 1999]

- Ideal cache model:
 - word size known to everyone in advance
 - two-level memory
 - Memory partitioned in lines of L words each
 - Cache size is Z
 - *Tall Cache Assumption*: $Z = \Omega(L^2)$
 - Eviction algorithm: *farthest in the future*¹
- *Cache-Oblivious Algorithm*: doesn't know Z, L
- Whenever it fetches word from disk, disk will send the line containing it (like in real world)

¹Can be replaced by LRU efficiently.

Why Cache-Oblivious Algorithms?

- Why should we care about this model?
 - Simpler code
 - Algorithm automatically “tunes” itself during execution
 - It is cool
 - Works for the many levels of memory hierarchy (all with their own parameters of Z and L)

Example: Median Finding

Algorithm: input Array of length n

① think of array partitioned into $\frac{n}{5}$ blocks of 5 elements each

② find median of each block (say by sorting)

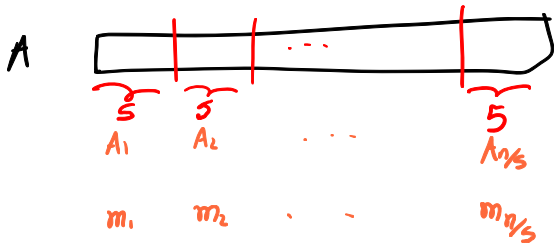
③ Find median of medians, call it x

④ Divide array into two $\begin{cases} \leq x \\ > x \end{cases}$

Recurse on array with most # elements.

Total work: $O(n)$

Memory used: ??



n

if $|C| > |D|$
 then median
 of A must be
 in C



$nir(B) = n/5$

$|C| \leq \frac{7n}{10}$
 elements

$MEDIAN(B) = x$
 $\geq 3n/10$

$\geq 3n/10$

A_{sorted}



$n/10$ m_i 's
 $\frac{1}{2} A_i$ is here

$n/10$ m_i 's
 $\frac{1}{2} A_i$ is here

Example: Median Finding $MT(n) \leftarrow$ # memory transfers that our algorithm will do on input size n

Algorithm, with memory handling:

- ① nothing to do (we are just thinking)
- ② find median of each block : can do this with linear scan through array $\therefore O(n/L + 1)$
- ③ find median of medians:
 - write each median to main memory $O(n/L + 1)$
 - recurse with our median finding algorithm $MT(n/5)$
- ④ Partition array into $\leq x, > x$
 - do this with a linear scan $O(n/L + 1)$
 - write arrays to memory $O(n/L + 1)$
- ⑤ Recurse $MT(7n/10)$

Example: Median Finding

Picture

Example: Median Finding - Recursion analysis

$$MT(n) = \underline{MT(n/5)} + \underline{MT(\frac{7n}{10})} + \underline{O(n/4 + 1)}$$

Base case: $MT(O(1)) = O(1)$ (i.e. ignoring cache)

$MT(n) \geq \underbrace{\# \text{ leaves in recursion}}_{L(n)}$

$$\underline{L(n)} = \underline{L(n/5)} + \underline{L(\frac{7n}{10})} \rightarrow \underline{n^\alpha = (n/5)^\alpha + (\frac{7n}{10})^\alpha}$$

$$\Rightarrow \alpha \approx 0.83$$

$$\therefore MT(n) \geq \underline{n^{0.83}} = \underline{\omega(n/4)} \text{ if } \underline{L = \omega(n^{0.17})}$$

Example: Median Finding - Recursion analysis

$$MT(n) = MT(n/5) + MT(\frac{7n}{16}) + O(n/L + 1)$$

Base case: $MT(O(L)) = O(1)$ (i.e. using cache)

\Rightarrow # leaves now is $L(n) = (n/2)^\alpha = O(n/L)$

Cost at each level decreases geometrically (practice problem)

\therefore $MT(n)$ dominated by root cost

$$\therefore MT(n) = O(n/L + 1)$$

- Administrivia
- Designing Algorithms in Real Life
- Cache-Aware and Cache-Oblivious Algorithms
- Conclusion
- Acknowledgements

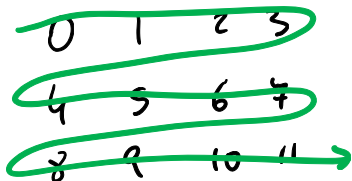
Cache-Aware Algorithm: Matrix Multiplication

Input : A, B $n \times n$ matrices

Output : $C = AB$

How to store matrices?

- major-row order: $A_{ij} > A_{i'j'}$ if $\begin{cases} i < i' \\ \text{or } i = i' \\ j < j' \end{cases}$



Cache-Aware Algorithm: Matrix Multiplication

- divide matrices A, B, C into blocks
(until they fit into cache)
then multiply blocks
- let s be size of the blocks
(parameter that must be set by algorithm)

$\text{BLOCK-MULT}_s(A, B, C, n)$

1 for $i \leftarrow 1$ to n/s

for $j \leftarrow 1$ to n/s

for $k \leftarrow 1$ to n/s

$\text{MULT}_s(A_{ik}, B_{kj}, C_{ij})$

Ordinary
Matrix
 $AB=C$
 $O(n^3)$ time

Cache-Aware Algorithm: Matrix Multiplication

- assume $s \mid n$ (otherwise more code to write)
- $s \times s$ matrices take $O(n^2/L)$ cache lines
- from full cache assumption ($Z = \Omega(L^2)$)

$s = O(\sqrt{Z})$ to minimize
cache complexity

$$s = \sqrt{Z/3} \quad \longmapsto \quad \underbrace{3 \times}_{A, B, C} \underbrace{\sqrt{\frac{Z}{3}} \cdot \sqrt{\frac{Z}{3}}}_{\text{size of matrix}} = Z$$

Cache-Aware Algorithm: Matrix Multiplication

- each call to MULT run with

$$O\left(\frac{z}{L}\right) = O\left(\frac{L^2}{z}\right) \text{ cache misses}$$

- $i, j, k \in \{1, \dots, \frac{n}{L}\}$ $\frac{n}{\sqrt{z}} = \frac{n}{L}$

algorithm has

$$O\left(\left(\frac{n}{\sqrt{z}}\right)^3 \cdot \frac{z}{L}\right) \text{ cache misses}$$

$$O\left(\frac{n^3}{L\sqrt{z}}\right)$$

Cache-Aware Algorithm: Matrix Multiplication

Cache-Aware Algorithm: Matrix Multiplication

Cache-Oblivious Algorithm: Matrix Multiplication

- use regular Matrix algorithm
- total work: $O(n^3)$
- matrices are given in major-row

Remark: using Strassen's algorithm $3 \rightarrow \log_2 7$

$$C = \begin{array}{|c|c|} \hline C_{11} & C_{12} \\ \hline C_{21} & C_{22} \\ \hline \end{array}$$

$$A = \begin{array}{|c|c|} \hline A_{11} & A_{12} \\ \hline A_{21} & A_{22} \\ \hline \end{array}$$

$$B = \begin{array}{|c|c|} \hline B_{11} & B_{12} \\ \hline B_{21} & B_{22} \\ \hline \end{array}$$

↓ ↓ ↓

$$C_{11} = A_{11} B_{11} + A_{12} B_{21}$$

$$C_{12} = A_{11} B_{12} + A_{12} B_{22}$$

..

$$C_{11} \leftarrow C_{11} + A_{11} B_{11}$$

if A_{11}, B_{11}, C_{11} don't fit in cache we recurse.

Cache-Oblivious Algorithm: Matrix Multiplication

$$MT(n) \leq 8 MT(n/2) + \underbrace{O\left(\frac{n^2}{L}\right)}_{\substack{\# \text{ cache misses} \\ \text{to put everything together}}}$$

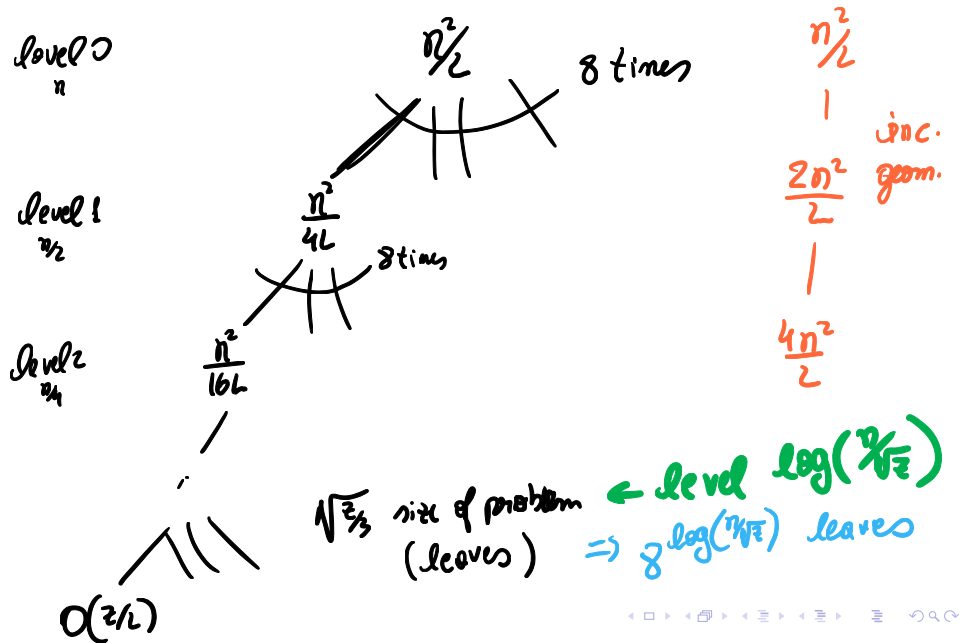
$$O\left(\frac{n^2}{L}\right)$$

linear scan
to add the
multiplied blocks

Base case: $MT\left(\sqrt{\frac{2}{3}}\right) = O\left(\frac{2}{L}\right)$

$\underbrace{\hspace{1cm}}$
cache misses
to bring matrices
to cache

Cache-Oblivious Algorithm: Matrix Multiplication



Cache-Oblivious Algorithm: Matrix Multiplication

$$O\left(\underbrace{8^{\log\left(\frac{n}{\sqrt{L}}\right)}}_{\# \text{ leaves}} \cdot \underbrace{\frac{L}{2}}_{\# \text{ cache misses per leaf}}\right) = O\left(\left(\frac{n}{\sqrt{L}}\right)^3 \cdot \frac{L}{2}\right)$$

$$= O\left(\frac{n^3}{L \cdot \sqrt{L}}\right).$$

Cache-Oblivious Algorithm: Matrix Multiplication

Cache-Oblivious Algorithm: Matrix Multiplication

Cache-Oblivious Algorithm: Matrix Multiplication

Conclusion

- We saw today how memory management affects execution of our algorithms

Conclusion

- We saw today how memory management affects execution of our algorithms
- Cache-Oblivious model: design algorithms with asymptotically optimal use of memory.
 - Simpler code
 - Efficient use of memory
 - Algorithm automatically “tunes” itself during execution
 - It is cool
 - Works for the many levels of memory hierarchy (all with their own parameters of Z and L)

Conclusion

- We saw today how memory management affects execution of our algorithms
- Cache-Oblivious model: design algorithms with asymptotically optimal use of memory.
 - Simpler code
 - Efficient use of memory
 - Algorithm automatically “tunes” itself during execution
 - It is cool
 - Works for the many levels of memory hierarchy (all with their own parameters of Z and L)
- Matrix Multiplication in cache-oblivious format

Acknowledgement

- Lecture based largely on:
 - [Frigo et al 1999]
 - 6.046 lecture notes

`http://stellar.mit.edu/S/course/6/sp15/6.046J/courseMaterial/
topics/topic2/lectureNotes/L23_-_Cache_Oblivious_I/L23_-_
Cache_Oblivious_I.pdf`

References I



Aggarwal, A. and Vitter, J.S., 1988.

The input/output complexity of sorting and related problems.

[Communications of the ACM](#), 31(9), pp. 1116-1127.



Frigo, M., Leiserson, C.E., Prokop, H. and Ramachandran, S., 1999.

Cache-oblivious algorithms.

In [40th Annual Symposium on Foundations of Computer Science](#) (pp. 285-297).
IEEE.